

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 07-230362

(43)Date of publication of application : 29.08.1995

(51)Int.Cl.

G06F 3/06

G11B 20/18

G11B 20/18

G11B 20/18

(21)Application number : 06-323920

(71)Applicant : HITACHI LTD

(22)Date of filing : 30.11.1994

(72)Inventor : TSUNODA HITOSHI  
TAKAMOTO YOSHIFUMI  
KAGIMASA TOYOHICO

(30)Priority

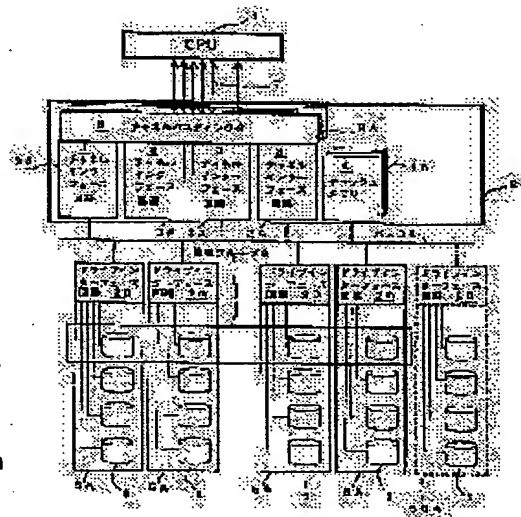
Priority number : 05329810 Priority date : 30.11.1993 Priority country : JP

## (54) DISK ARRAY DEVICE

(57)Abstract:

PURPOSE: To access data in a faulty drive by providing a circuit which writes data, which should be written in the faulty drive, in an alternative drive and writes data in normal drives on a board at the time when data write to the alternative drive is possible after substitution with the alternative drive of the faulty drive.

CONSTITUTION: If a fault occurs in any drive mounted on a drive board 5A, the faulty drive is not separated from the drive board 5A left connected to a mother board, but the drive board 5A is separated from the mother board with normal drives mounted. After the faulty drive is substituted with a normal alternative drive, the drive board 5A is connected to the mother board again. A logical group 9 is the fault recovery unit, and drives in the logical group 9 hold an error correction data group which consists of  $m-1$  data and a generated parity data.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-230362

(43) 公開日 平成7年(1995)8月29日

(51) Int.Cl. <sup>8</sup>	識別記号	片内整理番号	F I	技術表示箇所
G 0 6 F 3/06	5 4 0			
G 1 1 B 20/18	5 3 2 B	9074-5D		
	5 5 2 A	9074-5D		
	5 7 0 Z	9074-5D		

審査請求 未請求 請求項の数35 F D (全 25 頁)

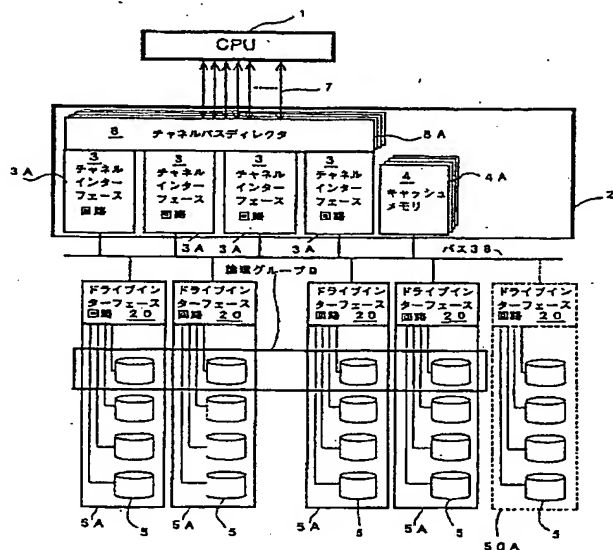
(21) 出願番号	特願平6-323920	(71) 出願人	000005108 株式会社日立製作所 東京都千代田区神田駿河台四丁目6番地
(22) 出願日	平成6年(1994)11月30日	(72) 発明者	角田 仁 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
(31) 優先権主張番号	特願平5-329810	(72) 発明者	高本 良史 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
(32) 優先日	平5(1993)11月30日	(72) 発明者	鍵政 豊彦 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
(33) 優先権主張国	日本 (J P)	(74) 代理人	弁理士 矢島 保夫

(54) 【発明の名称】 ディスクアレイ装置

(57) 【要約】 (修正有)

【目的】複数のドライブを基板に搭載し、これを、共通のマザーボードに並置でき、かつ何れかのドライブに障害が生じたときも、保持されたデータを上位装置からアクセス可能にするディスクアレイ装置を提供する。

【構成】アレイコントローラ2は、複数の基板の内、互いに異なるものに搭載されている複数のドライブに書き込む回路と、いずれかのドライブに障害が発生したためにマザーボードからその一つの基板が分離された状態において、その分離されたドライブに保持されているデータが読み出しのために上位装置からアクセスされたときに、他の複数の基板に保持されている複数のドライブから、その一つのデータと同じ誤り訂正データグループに属するデータおよびパリティデータを読み出す回路と、読み出された他の複数のデータと読み出されたパリティデータとから読み出すべきデータを回復する回路とを有する。



【特許請求の範囲】

【請求項 1】 (a) 一群のディスク記憶装置であって、該一群のディスク記憶装置は複数の論理グループに分割され、各論理グループは、複数のデータと、それらから生成された少なくとも一つの誤り訂正符号とをそれぞれ有する複数群の誤り訂正データを保持する複数のディスク記憶装置からなるものと、

(b) それぞれ該一群のディスク記憶装置の内の一部の複数のディスク記憶装置を保持する複数の個別基板と、

(c) 該一群の複数のディスク記憶装置の各々に分離可能に接続され、上位装置からのデータ読み出し要求またはデータ書き込み要求を実行するアレイコントローラであって、 (c 1) 各データ書き込み要求を、 (c 1 1) そのデータ書き込み要求に付随する書き込みデータが属する誤り訂正データ群に属すべき誤り訂正符号を生成し、 (c 1 2) 上記書き込みデータと該生成された誤り訂正符号とを、一つの誤り訂正データ群に属するデータとして、該複数の論理グループの一つに属する複数のディスク記憶装置に書き込むように実行し、 (c 2) 各データ読み出し要求を、 (c 2 1) それが要求するデータを保持する、該一群のディスク記憶装置の一つが正常なときには、そのディスク記憶装置からその要求されたデータを読み出し、 (c 2 2) そのディスク装置に障害があるときには、その障害があるディスク記憶装置が属する論理グループに属する複数のディスク記憶装置の内、上記障害があるディスク装置以外のディスク記憶装置に保持され、該読み出し要求で要求されたデータが属する、誤り訂正データ群に属し、一つの誤り訂正符号と他の複数のデータとから、該要求されたデータを回復するものとを有し、

(d) 該一群のディスク記憶装置の内、同一の論理グループに属する複数のディスク記憶装置が、該複数の個別基板の内の互いに異なるものに分散して搭載されているディスクアレイ装置。

【請求項 2】 該複数の個別基板を着脱可能に保持する、上記複数の個別の基板に共通な基板をさらに有し、上記アレイコントローラは、該一群の複数のディスク記憶装置に該共通の基板の上の信号線路を介して接続され、

該複数の個別基板は、該共通の基板にほぼ垂直に、かつ、互いにほぼ平行に保持されている請求項 1 記載のディスクアレイ装置。

【請求項 3】 各個別基板に搭載されている複数のディスク記憶装置は、互いに異なる複数の論理グループに属する請求項 1 記載のディスクアレイ装置。

【請求項 4】 各個別基板に搭載されている複数のディスク記憶装置は、互いに異なる複数の論理グループに属する請求項 2 記載のディスクアレイ装置。

【請求項 5】 上記アレイコントローラは、該障害がある一つのディスク記憶装置を搭載した上記複

数の個別基板の一つが、該共通の基板から分離されている後に上記上位装置から供給され、かつ、その障害があるディスク記憶装置に保持された一つのデータを要求する一つの読み出し要求を受理し、

50 その受理した読み出し要求の実行にあつては、その障害があるディスク記憶装置が属する論理グループに属し、その一つの基板以外の他の複数の基板に保持された他の複数のディスク記憶装置に保持された他の複数のデータからその要求されたデータを回復し、回復されたデータを上記上位装置に供給し、

さらに、上記分離の後に該上位装置から供給され、かつ、該一つの基板に保持された他の正常なディスク記憶装置に保持されたデータを要求する他の読み出し要求を受理し、

15 その受理した他の読み出し要求の実行にあつては、その正常な他のディスク記憶装置が属する他の論理グループに属し、その一つの基板以外の他の複数の基板に保持されたさらに他の複数のディスク記憶装置に保持されたさらに他の複数のデータからその要求された他のデータを回復し、回復された他のデータを上記上位装置に供給する請求項 2 記載のディスクアレイ装置。

【請求項 6】 ランダムアクセス可能なメモリをさらに有し、

該アレイコントローラは、

25 該分離の後に上記上位装置から供給され、該分離された一つの個別基板に保持されている複数の正常なディスク記憶装置もしくは該障害があるディスク記憶装置のいずれか一つへデータを書き込むことを要求する複数の書き込みを受理し、

30 該書き込みデータを該メモリに書き込み、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該メモリに保持された複数の書き込みデータを用いて該受理された複数の書き込み要求を実行する請求項 5 記載のディスクアレイ装置。

35 【請求項 7】 上記アレイコントローラは、該一つの個別基板が該共通に基板に再度接続された後に、上記ランダムアクセスメモリに保持された複数の書き込みデータの書き込みの前に、該一つの個別基板以外の該他の複数の基板に搭載された他の複数のディスク記憶装置内の複数のデータおよび複数の誤り訂正符号から、該障害があるディスク記憶装置に保持されていた複数のデータを回復し、該回復された複数のデータを該交替のディスク記憶装置に書き込む請求項 6 記載のディスクアレイ装置。

【請求項 8】 該アレイコントローラは、

45 該一つの個別基板が該共通の基板から分離されている状態において、上記上位装置から供給された上記書き込み要求を受理し、

その一つの個別基板が再度その共通の基板に接続されるを待たないで、その受理した書き込み要求を実行し、

50 その実行にあつては、その書き込みデータに対して誤

り訂正符号を生成し、その障害があるディスク記憶装置もしくはその正常なディスク記憶装置の一つと同じ論理グループに属し、その一つの個別基板以外の他の一つの個別基板に保持された一つのディスク記憶装置に生成した誤り訂正符号を書き込み、

該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換され、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該一つの個別基板に搭載された該交替のディスク記憶装置および該複数の正常なディスク記憶装置の各々書き込まれるべき複数のデータを、その一つの個別基板以外の他の複数の個別基板に保持された他の複数のディスク記憶装置に保持された他の複数のデータから回復し、

回復された複数のデータを、その、各ディスク記憶装置に書き込む請求項 5 記載のディスクアレイ装置。

【請求項 9】該アレイコントローラは、該一つの個別基板に保持された複数の正常なディスク記憶装置の各々に関しては、上記一つの個別基板が該共通基板から分離された後に上記上位装置から供給された複数の書き込み要求が指定した複数の書き込みデータを回復し、

該一つの個別の基板が該共通の基板から分離される前に該各正常なディスク記憶装置にすでに保持されていた書き込みデータは回復しない請求項 8 記載のディスクアレイ装置。

【請求項 10】該アレイコントローラに接続された少なくとも一つの予備のディスク記憶装置をさらに有し、該アレイコントローラは、

該一つの個別基板が該共通の基板から分離された状態で、該上位装置から供給された、該分離された一つの個別基板に搭載された該障害があるディスク記憶装置もしくは複数の正常なディスク記憶装置に書き込むべきデータを、上記予備のディスク記憶装置に書き込み、該上位装置からその後供給された読み出し要求が指定するデータが該予備のディスク記憶装置にすでに書き込まれているとき、該予備のディスク記憶装置からその要求されたデータを読み出し、該上位装置に供給する請求項 5 のディスクアレイ装置。

【請求項 11】該アレイコントローラは、該一つの個別基板が該共通の基板から分離された状態で、該上位装置から供給された読み出し要求を実行したときに、その読み出し要求が、該分離された一つの個別基板に搭載された、該障害があるディスク記憶装置もしくは複数の正常なディスク記憶装置に保持されたデータの読み出しを要求するとき、その読み出し要求の実行時に回復したデータを該予備のディスク記憶装置に書き込み、

該読み出し要求が指定するデータの読み出し要求がその後該上位装置から要求されたときに、該予備のディスク記憶装置からその要求されたデータを読み出す請求項 1

0 記載のディスクアレイ装置。

【請求項 12】該アレイコントローラは、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、上記予備ディスク記憶装置に保持された複数の書き込みデータの内、その個別基板が該共通基板から分離されていた間に該上位装置から供給された、その個別の基板に搭載された複数の正常なディスク装置に書き込まれるべきであった複数のデータを、該複数の正常なディスク記憶装置に転送する請求項 10 記載のディスクアレイ装置。

【請求項 13】該アレイコントローラは、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該障害ドライブに障害が発生する前に該障害があるディスク記憶装置にすでに保持されていたデータおよび該障害があるディスク記憶装置に障害が発生した後に該上位装置から供給された、その障害があるディスク記憶装置に書き込むべきであった他の複数の書き込みデータとを、該一つの個別基板以外の複数の個別基板に搭載された複数のディスク記憶装置に記憶された複数のデータと複数の誤り訂正符号から回復し、回復された複数のデータを該交替のディスク記憶装置に書き込む請求項 12 記載のディスクアレイ装置。

【請求項 14】該アレイコントローラに接続された複数の予備のディスク記憶装置をさらに有し、該アレイコントローラは、該一つの個別基板が該共通の基板から分離されたときに、該一つの個別基板に保持されている複数のディスク記憶装置の各々に対応して、そのディスク記憶装置の代わりに使用するディスク記憶装置を該複数の予備のディスク記憶装置から選択し、該一つの個別基板が該共通の基板から分離された状態で、該上位装置から供給された、該分離された一つの個別基板に搭載された該複数のディスク記憶装置のいずれか一つに書き込みデータを書き込むことを要求する書き込み要求を受理し、

その受理した書き込み要求の実行にあつては、上記複数の予備のディスク記憶装置の内、その一つのディスク記憶装置記憶に対応して選択されたディスク記憶装置にその書き込みデータを書き込み、該上位装置からその後供給された読み出し要求が指定するデータが上記複数の予備のディスク記憶装置のいずれかに保持されているとき、その予備のディスク記憶装置からその要求されたデータを読み出し、該上位装置に供給する請求項 5 記載のディスクアレイ装置。

【請求項 15】該共通の基板に着脱自在に、かつ、該複数の個別の基板に略平行に保持された少なくとも一つの予備の個別の基板をさらに有し、

該複数の予備のディスク記憶装置は、該予備の基板に搭載され、該共通の基板の上記信号線を介して該アレイコントローラに接続されている請求項 14 記載のディスクアレイ装置。

【請求項 16】該アレイコントローラは、  
該一つの個別基板が該共通の基板から分離された状態で、該上位装置から供給された読み出し要求を実行したときに、その読み出し要求が、該分離された一つの個別基板に搭載された、該障害があるディスク記憶装置もしくは複数の正常なディスク記憶装置に保持されたデータの読み出しを要求するとき、その読み出し要求の実行により回復したデータを該予備のディスク記憶装置に書き込み、  
該読み出し要求が指定するデータの読み出し要求がその後該上位装置から要求されたときに、該予備のディスク記憶装置からその要求されたデータを読み出す請求項 14 記載のディスクアレイ装置。

【請求項 17】該アレイコントローラは、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換され、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、上記複数の予備ディスク記憶装置の内、該一つの個別基板に搭載された複数の正常なディスク記憶装置に対応して選択された複数の予備のディスク記憶装置に保持された複数の書き込みデータを、該複数の正常なディスク記憶装置に書き込む請求項 14 記載のディスクアレイ装置。

【請求項 18】該アレイコントローラは、  
該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換され、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該一つの個別基板が該共通基板から分離される前に該上位装置から供給された該障害があるディスク記憶装置に書き込むべきであった複数の書き込みデータおよび該一つの個別基板が該共通基板から分離される前に該上位装置から供給された該障害があるディスク記憶装置に書き込むべきであった他の複数の書き込みデータとを、該一つの個別基板以外の複数の個別基板に搭載された複数のディスク記憶装置に保持された、複数のデータおよび複数の誤り訂正符号から回復し、

回復された複数のデータを該交替のディスク記憶装置に書き込む請求項 14 記載のディスクアレイ装置。

【請求項 19】該アレイコントローラは、  
該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換され、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該障害のあるドライブが障害状態になる前に該障害があるドライブが保持していた複数のデータを回復し、  
回復された複数のデータを該交替のディスク記憶装置に書き込み、

その書き込み該障害があるディスク記憶装置に対して選

択された該一つの予備のディスク記憶装置に保持されている、該一つの個別基板が該共通基板から分離された後に該上位装置から供給された該障害があるディスク記憶装置に書き込むべきであった複数の書き込みデータを、  
05 該交替ディスク記憶装置に転送する請求項 14 記載のディスクアレイ装置。

【請求項 20】該アレイコントローラは、  
該障害のあるディスク記憶装置が現に障害になった後に、該一つの個別基板上の該障害のあるディスク記憶装置が現に障害になる前に、その障害ドライブが保持していた複数のデータを、該一つの個別基板以外の複数の個別基板に搭載された複数のディスク記憶装置に保持された複数のデータから回復し、

10 回復された複数のデータを該障害のあるディスク記憶装置に対応して選択された予備のディスク記憶装置に書き込み、

該回復されたデータの書き込みが完了後、該上位装置から供給される、該障害があるディスク記憶装置に対する書き込み要求および読み出し要求を、該一つの予備のディスク記憶装置に対して実行するように、該読み出し要求および該書き込み要求の実行する請求項 17 記載のディスクアレイ装置。

【請求項 21】該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該障害があるディスク記憶装置に対応して選択された該一つの予備のディスク記憶装置を、該障害があるディスク記憶装置に代わりのディスク記憶装置として使用し、

20 30 該交替のディスク記憶装置を新たな予備のディスク記憶装置として使用する請求項 20 記載のディスクアレイ装置。

【請求項 22】該共通の基板に着脱自在に、かつ、該複数の個別の基板に略平行に保持された少なくとも一つの予備の個別の基板をさらに有し、  
35 該複数の予備のディスク記憶装置は、該予備の基板に搭載され、該共通の基板の上記信号線を介して該アレイコントローラに接続されている請求項 21 記載のディスクアレイ装置。

【請求項 23】該アレイコントローラは、  
該一つの個別基板が該共通の基板から分離されたときに、該一つの個別基板上の該障害のあるディスク記憶装置が保持していた書き込みデータを、該一つの個別基板以外の複数の個別基板に搭載された複数のディスク記憶装置に保持された複数のデータおよび複数の誤り訂正符号から回復し、

45 回復された複数のデータを該障害のあるディスク記憶装置に対応して選択された予備のディスク記憶装置に書き込み、

50 該回復されたデータの書き込みが完了後、該上位装置か

ら供給される、該障害があるディスク記憶装置に保持されているデータに対する書き込み要求および読み出し要求を、該一つの予備のディスク記憶装置に対して実行するように、該読み出し要求および該書き込み要求の実行する請求項17記載のディスクアレイ装置。

【請求項24】該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該分離された一つの個別基板が、該共通の基板に再度接続されたときに、該障害があるディスク記憶装置に対応して選択された該一つの予備のディスク記憶装置を、該障害があるディスク記憶装置に代わりのディスク記憶装置として使用し、

該交替のディスク記憶装置を新たな予備のディスク記憶装置として使用する請求項20記載のディスクアレイ装置。

【請求項25】該共通の基板に着脱自在に、かつ、該複数の個別の基板に略平行に保持された少なくとも一つの予備の個別の基板をさらに有し、

該複数の予備のディスク記憶装置は、該予備の基板に搭載され、該共通の基板の上記信号線路を介して該アレイコントローラに接続されている請求項24記載のディスクアレイ装置。

【請求項26】該アレイコントローラは、該上位装置から供給された書き込み要求に付随する一つの書き込みデータを複数のサブデータに分割し、該生成された複数のサブデータに対する誤り訂正符号を生成し、該複数のサブデータと該生成された誤り訂正符号を一つの誤り訂正データ群に属するデータとして、互いに異なる個別基板に搭載された複数のディスク記憶装置に書き込むように、該書き込み要求を実行し、さらに、該上位装置から供給されたデータ読み出し要求が指定するデータを構成する複数のサブデータを、異なる個別基板に搭載された複数のディスク記憶装置から読み出し、読み出された複数のサブデータを結合し、結合後のデータを該上位装置に供給するように、各読み出し要求を実行する請求項4記載のディスクアレイ装置。

【請求項27】各誤り訂正データ群は、該上位装置から供給された複数の書き込みデータと、該複数の書き込みデータに対する誤り訂正符号とからなり、

該アレイコントローラは、上位装置から供給された各書き込み要求に実行にあつては、その書き込み要求に付随する書き込みデータが属する誤り訂正符号データ群に対してすでに決定された旧の誤り訂正符号と、該書き込みデータで更新される旧書き込みデータとを、それぞれ互いに異なる個別基板に保持された相異なるディスク記憶装置から読み出し、読み出された旧誤り訂正符号と読み出された旧書き込みデータと該書き込みデータとから、該書き込みデータをもって該旧書き込みデータを更新した後の誤り訂正符号を生成し、

その書き込みデータと、その生成された誤り訂正符号を互いに異なる個別基板に搭載された複数のディスク記憶装置に書き込み、

上位装置から供給された読み出し要求の実行にあつては、その読み出し要求が指定するデータを、いずれかひとつの個別基板に搭載された一つのディスク記憶装置から読み出し、読み出されたデータを該上位装置に供給する請求項4記載のディスクアレイ装置。

【請求項28】上記アレイコントローラは、該障害がある一つのディスク記憶装置を搭載した上記複数の個別基板の一つが、該共通の基板から分離されている後に上記上位装置から供給され、かつ、その障害があるディスク記憶装置に保持された一つのデータを要求する一つの読み出し要求を受理し、

その受理した読み出し要求の実行にあつては、その障害があるディスク記憶装置が属する論理グループに属し、その一つの基板以外の他の複数の基板に保持された他の複数のディスク記憶装置に保持された他の複数のデータからその要求されたデータを回復し、回復されたデータを上記上位装置に供給し、

さらに、上記分離の後に該上位装置から供給され、かつ、該一つの基板に保持された他の正常なディスク記憶装置に保持されたデータを要求する他の読み出し要求を受理し、

その受理した他の読み出し要求の実行にあつては、その正常な他のディスク記憶装置が属する他の論理グループに属し、その一つの基板以外の他の複数の基板に保持されたさらに他の複数のディスク記憶装置に保持されたさらに他の複数のデータからその要求された他のデータを回復し、回復された他のデータを上記上位装置に供給する請求項1記載のディスクアレイ装置。

【請求項29】該アレイコントローラは、該分離の後に上記上位装置から供給され、該分離された一つの個別基板に保持されている複数の正常なディスク記憶装置もしくは該障害があるディスク記憶装置のいずれか一つへデータを書き込むことを要求する複数の書き込み要求を受理し、

該書き込みデータを該メモリに書き込み、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後、該メモリに保持された複数の書き込みデータを用いて該受理された複数の書き込み要求を実行する請求項28記載のディスクアレイ装置。

【請求項30】該アレイコントローラに接続された複数の予備のディスク記憶装置をさらに有し、

上記アレイコントローラは、該一つの個別基板が該共通の基板から分離されたときに、該一つの個別基板に保持されている複数のディスク記憶装置の各々に対応して、そのディスク記憶装置の代わりに使用するディスク記憶装置を該複数の予備のディ



スク記憶装置から選択し、

該一つの個別基板が該共通の基板から分離された状態で、該上位装置から供給された、該分離された一つの個別基板に搭載された該複数のディスク記憶装置のいずれか一つに書き込みデータを書き込むことを要求する書き込み要求を受理し、

その受理した書き込み要求の実行にあつては、上記複数の予備のディスク記憶装置の内、その一つのディスク記憶装置記憶に対応して選択されたディスク記憶装置にその書き込みデータを書き込み、

該上位装置からその後供給された読み出し要求が指定するデータが上記複数の予備のディスク記憶装置のいずれかに保持されているとき、その予備のディスク記憶装置からその要求されたデータを読み出し、該上位装置に供給する請求項 28 記載のディスクアレイ装置。

【請求項 31】一群のディスク記憶装置であつて、該一群のディスク記憶装置は複数の論理グループに分割され、各論理グループは、複数のデータと、それらから生成された少なくとも一つの誤り訂正符号とをそれぞれ有する複数の誤り訂正データを保持する複数のディスク記憶装置からなるものと、

それぞれ該一群のディスク記憶装置の内一部の複数のディスク記憶装置を保持する複数の個別基板と、該一群の複数のディスク記憶装置の各々に分離可能に接続され、上位装置からのデータ読み出し要求またはデータ書き込み要求を実行するアレイコントローラとを有するディスクアレイ装置において、

(a) 該一群のディスク記憶装置の内、同一の論理グループに属する複数のディスク記憶装置は、該複数の個別基板の内の互いに異なるものに分散して搭載された複数のディスク記憶装置から構成されるように、各論理グループに属するディスク記憶装置を決定し、

(b) いずれかのディスク記憶装置に障害が発生したとき、そのディスク記憶装置を含む一つの個別基板を該アレイコントローラから分離し、

(c) 該アレイコントローラから分離されている状態で上記上位装置から発行され、該分離された一つの個別の基板上の上記障害が発生したディスク記憶装置および他の正常なディスク記憶装置の内のいずれか一つに保持されたデータを要求する複数の読み出し要求を受理し、

(d) 各受理した読み出し要求が要求する読み出しデータを、その一つの個別基板以外の基板に保持された複数のディスク記憶装置から回復し、回復されたデータを該上位装置に供給するように、各読み出し要求を実行するディスクアレイ装置のアクセス方法。

【請求項 32】該アレイコントローラはランダムアクセス可能なメモリをさらに有し、該方法は、該アレイコントローラから分離されている状態で上記上位装置から発行され、該分離された一つの個別基板に保持されている該障害があるディスク記憶装置もしくはは

複数の正常なディスク記憶装置のいずれか一つへ書き込むべきデータを指定する書き込み要求を受理し、該受理した書き込み要求が指定する書き込みデータを、該メモリに一時的に保持し、

05 該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後に、該一つの個別基板が該アレイコントローラに再度接続されたときに、該メモリに保持された該書き込みデータを用いて該一つの書き込み要求を実行するステップをさらに有する請求項 31 記載のディスクアレイ装置のアクセス方法。

【請求項 33】上記アレイディスク装置は、該アレイコントローラに接続された少なくとも一つの予備のディスク記憶装置をさらに有し、

上記方法は、

15 該一つの個別基板が該共通の基板から分離された状態で、該上位装置から受理され、該分離された一つの個別基板に搭載された該障害があるディスク記憶装置もしくは他の正常なディスク記憶装置のいずれか一つに書き込むべきデータを指定する複数の書き込み要求を受理し、

20 該書き込み要求の各々を、それが指定する書き込みデータを上記予備のディスク記憶装置に書き込むように実行し、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換され、該一つの個別基板が該アレイコントローラに再度接続された後に、該予備のディスク記憶装置に保持された書き込みデータの内、上記一つの個別基板に保持された該他の正常なディスク装置に書き込まれるべきデータを、該他の正常なディスク記憶装置に転送するステップをさらに有する請求項 31 記載のディスクアレイ装置のアクセス方法。

30 【請求項 34】上記アレイディスク装置は、該アレイコントローラに接続された複数の予備のディスク記憶装置をさらに有し、

35 該一つの個別基板の上の各ディスク装置に対応して、該複数の予備のディスク記憶装置の一つを選択し、

該一つの個別基板が該アレイコントローラから分離されるときに、該障害が発生した一つのディスク記憶装置が保持していたデータを、該一つの個別基板以外の他の個別基板に保持された複数のデータおよび複数の誤り訂正符号から回復し、

回復された複数のデータを該一つの障害が発生したディスク記憶装置に対応して選択された、一つの予備のディスク記憶装置に転送し、

45 該一つの個別基板が該アレイコントローラから分離された状態で、該上位装置から受理したいずれかの書き込み要求が、該分離された一つの個別基板に搭載された該障害があるディスク記憶装置もしくは複数の正常なディスク記憶装置のいずれか一つに書き込むべきデータを指定するとき、上記複数の予備のディスク記憶装置の内、そ

の一つのディスク記憶装置に対応して選択された一つの予備のディスク記憶装置に書き込み、該障害があるディスク記憶装置が正常な交替のディスク記憶装置により置換された後に、該一つの個別基板が該アレイコントローラに再度接続されたときに、該一つの予備のディスク記憶装置を、該障害が発生した一つのディスク記憶装置の代わりに使用し、該交替のディスク記憶装置を新たな予備のディスク記憶装置として使用するステップをさらに有する請求項 3 1 記載のディスクアレイ装置のアクセス方法。

【請求項 3 5】該一つの個別基板が該アレイコントローラに再度接続された後に、該複数の予備のディスク記憶装置の内、該他の正常なディスク記憶装置に対して選択された他の一つの予備のディスク記憶装置に保持された書き込みデータを、該他の正常なディスク記憶装置に転送し、

上記他の予備のディスク記憶装置を新たに予備のディスク記憶装置として解放するステップをさらに有する請求項 3 4 記載のディスクアレイ装置のアクセス方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、複数のデータを誤り訂正符号とともに保持する複数のディスク装置を有するディスクアレイ装置に関する。

【0002】

【従来の技術】現在のコンピュータシステムにおいては、CPU 等の上位側が必要とするデータは 2 次記憶装置に格納され、CPU が必要とするときに 2 次記憶装置に対してデータの書き込み、読み出しを行うようになっている。この 2 次記憶装置としては、一般に不揮発な記憶媒体が使用され、代表的なものとして磁気ディスク装置や光ディスクなどがあげられる。以下、これらに例示されるディスク装置をドライブと呼ぶ。

【0003】近年の高度情報化に伴い、コンピュータシステムにおいて、2 次記憶装置の高性能化が要求されてきた。その一つの解として、多数の比較的容量の小さなドライブにより構成されるディスクアレイが考えられている。

【0004】例えば、D. Patterson, G. Gibson, and R. H. Kartz: A Case for Redundant Array of Inexpensive Disks (RAID), ACM SIGMOD Conference, Chicago, IL, (June 1988), pp. 109-116 (以下、第 1 の従来技術と呼ぶ) なる論文においては、データを分割して並列に処理を行うディスクアレイ (レベル 3) とデータを分散して独立に扱うディスクアレイ (レベル 4, 5) について、それらの性能および信頼性の検討結果が報告されている。

【0005】レベル 3 のディスクアレイは、CPU から一つの書き込み要求に付随して転送されてきた一つの書き込みデータを複数の分割し、分割後の複数のデータから誤り訂正用のパリティデータを作成する。そして、こ

れらのデータ (分割後の複数のデータおよびパリティデータ) を複数のドライブに並列に格納する。また、CPU からデータの読み出しが指示された場合は、まず各々のドライブから分割されたデータを並列に読み込み、それらを結合して CPU へ転送する。

【0006】レベル 4, 5 のディスクアレイは、CPU から転送されてきた複数の書き込み要求に付随する複数のデータを上記分割したデータの代わりに使用すること、およびデータの読み出しは、個々のデータに対して行なわれ、レベル 3 でのデータの結合は行なわれない点で、レベル 3 の場合と異なる。

【0007】複数のデータと、それらから生成された誤り訂正用のデータとは、「誤り訂正データ群」あるいは「パリティグループ」あるいは「パリティデータグループ」と呼ばれることがある。誤り訂正用のデータとしてはパリティデータのほかにも各種の誤り訂正用データを用いることができるが、本明細書でも、「パリティグループ」との用語を誤り訂正用データがパリティデータでない場合にも用いるものとする。また、以下では誤り訂正符号としては、パリティデータの場合について説明するが、本願発明はこのパリティデータ以外にも適用可能である。

【0008】以下では、同一のパリティグループのデータを格納する複数のドライブを「論理グループ」と呼ぶ。論理グループは障害回復の単位であり、論理グループ内の何れかのドライブに障害が発生した場合はその論理グループ内の他のドライブのデータから障害回復を行うことができる。

【0009】例えば、分割後のデータを格納したドライブの中の 1 台に障害が発生し、データが読み出せなくなった場合は、残りのドライブ内のデータとパリティデータとから、障害が発生したドライブ内のデータを復元することができる。ディスクアレイのような多数のドライブにより構成される装置では、部品点数が増加することにより、障害が発生する確率が高くなる。そこで、ディスク装置としての信頼性の向上を図る目的で、このようにパリティを用意して障害の回復を行うことができるようにしている。

【0010】また、何れかの障害ドライブを他の交替ドライブと置換するまでの間、その障害ドライブに書き込むべきデータを、あらかじめ設けられた予備のドライブに書き込むことも提案されている。例えば、米国特許 5, 077, 736 (以下、第 2 の従来技術と呼ぶことがある) に開示されたものである。

【0011】従来の、典型的なレベル 3 のディスクアレイ装置 (以下、第 3 の従来技術と呼ぶことがある) では、複数のドライブは、ケーブルでもってアレイコントローラに接続される。その際、異なる論理グループの各々を構成するドライブの数に等しい数のバスでもってそれらのドライブをアレイコントローラに接続していた。



すなわち、異なる論理グループの第1のドライブが、共通の第1のバスにてアレイコントローラに接続されるように、それらのドライブはケーブルでもってデータチェーン状に接続される。同様に、それらの論理グループの第2のドライブが、共通の第2のバスにてアレイコントローラに接続されるように、それらのドライブはケーブルでもってデータチェーン状に接続される。他のドライブも同様である。このようなケーブルを使用したドライブの接続は、装置を構成するドライブの総数が多くなるときには、装置専有面積あるいは保守の容易さなどの点であまり望ましくない。したがって、小型のドライブを多数実装する方法が望まれる。

【0012】特開平3-108178号(米国特許出願番号第409495号(1989年9月19日出願)に対応する)(以下、第4の従来技術と呼ぶことがある)には、底面に多数のピンを有する複数のドライブを半導体基板に搭載する技術が開示されている。この技術は、小型のドライブを多数基板上に実装することを目的としている。しかし、この従来技術には、これらのドライブにより論理グループを構成するか否か、構成するときにはどのようなドライブにより各論理グループを構成するかについては開示されていない。

【0013】特開平4-228153号(米国特許出願番号第502215号(1990年3月30日出願)に対応する)(以下、第5の従来技術と呼ぶことがある)には、それぞれ複数のドライブを共通の配線付きの基板の一つの面に着脱可能に搭載した複数の装置(アレイと呼ばれている)を使用したディスクアレイが開示されている。各アレイ内の複数のドライブは、同一の論理グループに属する。何れかの基板の何れかのドライブに障害が生じたとき、そのドライブを搭載した基板をシステムに接続したままに保持し、その障害が生じたドライブを他の交替ドライブと置換するためにその基板から分離する。しかし、この状態でも、その障害ドライブへのアクセスが可能になっている。すなわち、その障害ドライブに保持されているデータの読み出し要求が発生したとき、その論理グループ内の他のドライブ内のデータとパリティデータとから、そのアクセス対象となったデータを回復して読み出す。

【0014】

【発明が解決しようとする課題】上記第5の従来技術では、複数のドライブを共通の基板に搭載したアレイを共通のマザーボードに搭載するので、多くの小型のドライブを高密度に実装することが可能になる。しかし、何れかのドライブに障害が発生したときに、その障害ドライブを他の交替ドライブと置換するために、その障害ドライブがそれを搭載した基板から分離できなければならない。この障害ドライブの分離のためには、各ドライブの上方に、ユーザが手でそのドライブをアクセスできる空間を残す必要がある。したがって、この空間のために、

このような基板を狭い間隔で多数配列するには限度がある。

【0015】したがって、本発明の目的は、複数のドライブを複数の基板に搭載し、このような複数の基板を、共通のマザーボードに比較的狭い間隔で並置でき、それでいて何れかのドライブに障害が生じたときでも、その障害ドライブに保持されたデータを上位装置からアクセス可能にするディスクアレイ装置を提供することにある。

【0016】

【課題を解決するための手段】この目的達成のために、本発明によるディスクアレイ装置は、マザーボードと、該マザーボードの一つの面に着脱可能に保持され、互いにはほぼ平行に配置された複数の基板と、各基板の一つの面に搭載された複数のドライブと、該マザーボード上に設けられた信号線路を介して該複数の基板に搭載された複数のドライブからのデータの読み出しおよびそれらへのデータの書き込みを制御するアレイコントローラとを有し、該アレイコントローラは、同一の誤り訂正データグループに属する、複数のデータおよびそれらの複数のデータに対して使用する誤り訂正コードを、上記複数の基板の内、互いに異なるものに搭載されている複数のドライブに書き込む回路と、上記基板のいずれか一つの上に搭載されたいずれかのドライブに障害が発生したために該マザーボードからその一つの基板が分離された状態において、その分離された一つの基板の上のドライブに保持されているデータが読み出しのために上位装置からアクセスされたときに、そのドライブが障害ドライブの場合でもあるいは正常なドライブでも、それぞれその一つの基板以外の他の複数の基板に保持されている複数のドライブから、その一つのデータと同じ誤り訂正データグループに属する他の複数のデータおよびパリティデータを読み出す回路と、読み出された他の複数のデータと読み出されたパリティデータとから該読み出すべき一つのデータを回復する回路とを有する。

【0017】本願発明の第1の望ましい態様では、上記障害ドライブを搭載した一つの基板がマザーボードから分離されている状態で、その一つの基板に搭載されたいずれか一つのドライブに書き込むべきデータが上位装置から供給されたときに、そのデータを一時的に保持するランダムアクセス可能なメモリと、上記障害ドライブがいずれかの交替ドライブにより置換された後、上記一つの基板が上記マザーボードに再度接続され、その交替ドライブがデータ書き込み可能になったときに、上記メモリに保持された書き込みデータの内、上記障害ドライブに書き込むべきデータを該交替ドライブに書き込み、該分離されていた基板上の正常なドライブに書き込むべきデータを、その正常なドライブに書き込む回路を有する。

【0018】本発明の他の望ましい態様では、予め上記

データを保持するためのドライブが搭載されている基板と異なる他の基板上に設けられた、少なくとも一つの基板上の搭載されるデータ保持用のドライブの数以上の予備のドライブと、上記障害ドライブを搭載した一つの基板がマザーボードから分離されたとき、その一つの基板上の各ドライブに対応して、該複数の予備のドライブの一つを選択し、その一つの基板が分離された状態で、その一つの基板に搭載されたいずれか一つのドライブに書き込むべきデータが上位装置から供給されたときに、上記複数の予備のドライブの内、そのドライブに対して選択された予備のドライブにその書き込みデータを書き込む回路と、上記障害ドライブがいずれかの交替ドライブにより置換された後、上記一つの基板が上記マザーボードに再度接続され、その交替ドライブがデータ書き込み可能になったときに、該一つの基板の上の各ドライブに上記複数の予備のドライブの内、その各ドライブに対して選択された予備のドライブに保持された書き込みデータを書き込む回路とが設けられる。

【0019】

【作用】 いずれかのドライブに障害が発生し、この障害が発生したドライブを正常なドライブに交換する際に、基板をマザーボードのコネクタから外しても、他の基板内のドライブのデータを用いて、取り外した基板のドライブのデータを回復できる。したがって、障害ドライブを含む基板をマザーボードのコネクタから外し、障害ドライブを正常ドライブに交換する作業を行っている間であっても、CPUからの読み出しや書き込み処理は中断せずにすむ。

【0020】

【実施例】 以下、図面を用いて本発明の実施例を説明する。

【0021】 (実施例1)

(1) 装置の概要

図1は、本発明に係るディスクアレイシステムの第1の実施例の全体構成を示す。この図において、ディスクアレイシステムは、複数のドライブ5とそれらの動作を制御するアレイコントローラ2とからなる。

【0022】 1は外部からデータの書き込みおよび読み出しを指令するCPUである。CPU1とアレイコントローラ2とは、複数本のチャネルバス7により接続される。アレイコントローラ2には、チャネルバス7を選択するチャネルバスディレクタ8を分散して搭載した、複数、例えば4つのチャネルバスディレクタボード8Aと、CPU1とのデータ転送制御を行うためのチャネルインターフェース回路3をそれぞれ搭載した複数、例えば4つのチャネルバスインターフェースボード3Aと、バッテリバックアップ等により不揮発化された半導体メモリであるキャッシュメモリ4を分散して保持する複数、例えば4つのキャッシュメモリボード4Aと、これらのキャッシュメモリボード4Aとチャネルインターフ

ェースボード3Aとを接続するバス38とが設けられている。このバス38に接続して、それぞれ複数のドライブ5とドライブインターフェース回路20をそれぞれ保持する複数、例えば、4つのドライブボード5Aが設けられている。

【0023】 図2(a)は、図1の装置に使用される複数のボードの配置を示す図である。各種のボードは、図2(a)に示すように、共通のマザーボード37に搭載されている。図2(a)において、各ボード上に図示された四角形370は、各ボードに搭載されたLSIチップを示す。ただし、ドライブ用のボード5Aの上に図示された四角形の一部370はLSIチップであるが残りの大部分5は、そのボードに搭載されたドライブを示す。このマザーボード37上には、図2(b)に示す、各ボード用のコネクタ群30A、40A、50A、80Aが設けられ、各ボードはこのようなコネクタによりマザーボードに対して着脱自在に保持されている。さらに、各ボードは互いにほぼ平行に配置されている。

【0024】 さらにマザーボード表面には、チャネルインターフェース用コネクタ群30A、キャッシュメモリ用コネクタ群40A、ドライブボード用のコネクタ群50Aを接続するように、前述のバス38が形成されている。さらに、チャネルディレクタ用のコネクタ群80Aとチャネルインターフェース用のコネクタ群30Aとを接続する信号配線12(図3)も設けられている。

【0025】 本実施例の特徴は、何れかのドライブに障害が発生したとき、そのドライブを搭載したドライブボードをマザーボード37に接続したまま障害ドライブをそのドライブボードから分離するのではなく、そのドライブボードを、正常なドライブを搭載したままマザーボード37から分離し、その障害ドライブを正常な交替ドライブにより置換した後、そのドライブボードをマザーボード37に再度接続することを可能にした点にある。

【0026】 したがって、何れかのドライブに障害が発生したときに、そのドライブを搭載したドライブボードを抜き差しすればよいので、各ドライブボード間の間隔を小さくすることができる。しかも、複数のドライブボードに分散して搭載された複数のドライブでもって一つの論理グループが構成されている。これにより、何れかのドライブボードがマザーボードから分離された状態でも、その分離されたドライブボード上に搭載されたドライブに保持されたデータを、他のドライブボードに保持された複数のドライブに保持されたデータおよびパリティデータから回復可能になっている。

【0027】 さて、論理グループ9は、複数(本実施例では4枚)のボード間に渡るm台(本実施例では4台)のドライブ5により構成される。論理グループ9は、障害回復単位であり、この論理グループ9内のドライブ5は、m-1個のデータとそれらから生成したパリティデータとからなる誤り訂正データグループ(すなわち、パ

リティグループ)を保持する。

【0028】なお、本実施例では、チャンネルインターフェースボード3A、ドライブボード5Aの数を各々4枚としたが、本発明はこれに限定されるものではない。また、ボードに装着されるドライブ5の数、および論理グループ9を構成するドライブ5の数も、本実施例に限定されず、任意の数としてよい。

【0029】図3は、図1の各ボード内の回路の概略構成を示す。

【0030】チャンネルバスディレクタ8は、各チャンネルバス7に接続されたインターフェースアダプタ(IF Adp)10、およびチャンネルバススイッチ11を備えている。このチャンネルバスディレクタ8は、実際は4つのボードに分散して配置されるが、図では簡単化のために、この回路の全体を示してある。

【0031】チャンネルインタフェース回路3は、CPU1から転送されてきたデータに対しプロトコル変換と転送速度を調整するチャンネルインターフェース(CH IF)13、データの転送制御を行なうデータ制御回路(DCC)14、CPU1から送られたデータをキャッシュメモリ4へ書き込むのを制御するチャンネル側キャッシュアダプタ(C Adp)15、このデータを分割して得られる複数のサブデータからパリティを生成するパリティ生成回路18、およびこれらの回路を制御するマイクロプロセッサ(MP)17を備えている。チャンネルインタフェース13は、データ線12によりチャンネルバススイッチ11と接続されている。キャッシュアダプタ15は、バス38の制御と、CPU1、キャッシュ4間のデータ転送制御を行なう回路である。

【0032】ドライブボード5Aには、ドライブ側キャッシュアダプタ(C Adp)19、およびドライブインターフェース回路(Drive IF)20が設けられている。16は、チャンネルバススイッチ11、チャンネルインタフェース13、データ制御回路14、チャンネル側キャッシュアダプタ15、マイクロプロセッサ17、キャッシュメモリ4、ドライブ側キャッシュアダプタ19、およびドライブインターフェース回路20を接続する制御信号線である。バス38は、制御情報を転送する制御バスと、データを転送するデータバスからなる。

【0033】なお、図1～図3においてボード50Aは後の実施例2で使用するもので、予備のドライブ5を搭載したボードである。

【0034】(2) 障害ドライブがないときの装置動作の概要

次に、図3を参照して、ディスクアレイシステムの内部動作を説明する。このときの動作は、基本的には、RAID3の原理にしたがった動作である。

【0035】CPU1より発行された読み出しまたは書き込みコマンドは、チャンネルバス7を通過して、アレイコントローラ2のチャンネルバスディレクタ8に入力する。

【0036】CPU1からコマンドが発行されると、アレイコントローラ2内のチャンネルバスディレクタ8によりコマンドの受付が可能かどうか判断する。CPU1からアレイコントローラ2に送られてきたコマンドはインターフェースアダプタ(IF Adp)10により取り込まれ、マイクロプロセッサ(MP)17はコマンドの受け付け処理を行なう。本実施例では、基本的に、従来のレベル3のRAIDによる方法にしたがってデータの書き込みおよび読み出しを行う。

【0037】図6は、RAIDレベル3によるデータの書き込みの様子を示す説明図である。図6において、BADR<sub>i</sub> (i=1, 2, 3または4)は、それぞれ、4つのドライブボード5Aのボードアドレスを示す。SD<sub>#i</sub> (i=1, 2, 3または4)は、各ボード内の1つのドライブを特定するドライブアドレスを示す。

【0038】以下、説明の便宜のため、例えばボードアドレスがBADR1であるボードは、ボードBADR1と呼ぶ。また、ボード内のドライブアドレスがSD#1のドライブは、ドライブSD#1と呼ぶ。

【0039】ボードBADR1のドライブSD#1、ボードBADR2のドライブSD#1、ボードBADR3のドライブSD#1、およびボードBADR4のドライブSD#1により、論理グループが構成されているものとする。また、ボードBADR4のドライブは、パリティデータの格納に用いるものとし、他のボードBADR1からBADR3内のドライブSD#1は、データの格納に用いるものとする。

【0040】CPU1からのコマンドがデータの書き込みコマンドであり、CPU1から転送された書き込みデータが、図6のD#1であるものとする。D#1は、その書き込みデータを特定する論理アドレスである。以下、論理アドレスD#1で特定されるデータをデータD#1と呼ぶ。

【0041】本実施例では、CPU1から書き込みまたは読み出すデータのデータ長は常に4KBの大きさの固定長データであるものとする。したがって、CPU1からの書き込みデータD#1も4KBである。なお、このデータ長により本発明が限定されることは無い。

【0042】CPU1からの書き込みデータD#1は、図1のキャッシュメモリ4に保持される。図6に示すように、RAID3ではこの保持された書き込みデータD#1を、データの先頭から1バイトずつ、論理グループを構成する4台のドライブのうちデータ格納用の3台のドライブ(BADR1, BADR2, BADR3のそれぞれのSD#1)に振り分けて分割する。分割したデータをサブデータと呼び、それぞれD#1-1, D#1-2, D#1-3というサブデータ名を付ける。なお、サブデータ名がD#1-1であるサブデータを、サブデータD#1-1と呼ぶ。

【0043】図6から分かるように、サブデータD#1

ー1はデータD#1の第1バイト目、第4バイト目、第7バイト目、…からなり、サブデータD#1-2はデータD#1の第2バイト目、第5バイト目、第8バイト目、…からなり、サブデータD#1-3はデータD#1の第3バイト目、第6バイト目、第9バイト目、…からなる。

【0044】サブデータに分割した後、これらのサブデータからパリティを作成する。このパリティは、サブデータにおいて各々対応するバイトに対して作成する。すなわち、サブデータD#1-1の第1番目のバイトとD#1-2の第1番目のバイトとD#1-3の第1番目のバイトとからパリティを求め、サブデータD#1-1の第2番目のバイトとD#1-2の第2番目のバイトとD#1-3の第2番目のバイトとからパリティを求め、…というようにパリティを生成する。生成したパリティを順に並べたパリティデータをP#1と呼ぶ。

【0045】これらのサブデータD#1-1、D#1-2、D#1-3およびパリティデータP#1を、論理グループを構成する4台のドライブに、分散して書き込む。ここでは、ボードBADR1のドライブSD#1にサブデータD#1-1を書き込み、ボードBADR2のドライブSD#1にサブデータD#1-2を書き込み、ボードBADR3のドライブSD#1にサブデータD#1-3を書き込み、ボードBADR4のドライブSD#1にパリティデータP#1を書き込むこととなる。

【0046】以上のようにして、CPU1からのデータの書き込みが行われる。

【0047】CPU1からデータの読み出しが要求されたときには、その要求されたデータを保持する論理グループから、そのデータを構成するサブデータを読み出し、それらのサブデータを結合して、キャッシュメモリ4を介してCPU1に送る。

【0048】(3) アドレス変換用テーブル

以上のデータの書き込みおよび読み出しに必要なアドレス変換テーブルの詳細を以下に説明する。

【0049】本実施例では、ディスクアレイシステムを構成するドライブ5としてSCSIインターフェースのドライブを使用する。また、本実施例では、CPU1はディスクアレイを意識せずに書き込みコマンドや読み出しコマンドを発行するものとする。すなわち、CPU1は、書き込みコマンドや読み出しコマンドを発行する相手であるディスク装置が従来方式のものかディスクアレイ装置であるかは意識していない。これは、現在主流のディスク装置が従来方式のものであり、OS(オペレーティングシステム)はディスクアレイを意識したものになっていないからである。

【0050】このため、本実施例では、CPU1では従来方式のドライブに対しデータを読み出しまたは書き込みを行っているものとして処理を行い、ディスクアレイシステムにおいて独自にディスクアレイで処理を行って

いる。したがって、CPU1は従来のインターフェースで書き込みや読み出しのコマンドを発行する。

【0051】具体的には、CPU1からアレイコントローラ2にデータの書き込みを指示する場合、CPU1は、書き込みデータに論理アドレス(データ名またはデータ番号)を付して転送する。アレイコントローラ2では、図6で説明したように、この書き込みデータをサブデータに分割し、パリティデータを作成し、それらのサブデータおよびパリティデータを論理グループ9の各ドライブに格納する。この際、それらのサブデータおよびパリティデータには、アレイコントローラ2において処理するアドレス(図6で説明したD#1-1~D#1-3やP#1など)がつけられる。CPU1から与えられる論理アドレスとこのアレイコントローラ2におけるアドレスとの対応をとるテーブルが、以下で説明するアドレス変換用テーブルである。

【0052】CPU1からアレイコントローラ2にデータの読み出しを指示する場合、CPU1は、論理アドレスを与えてアレイコントローラ2に読み出しを指示する。アレイコントローラ2は、以下で説明するアドレス変換用テーブルを利用して論理アドレスをディスクアレイシステム内のアドレスに変換し、そのアドレスでサブデータを読み出して結合してCPU1に渡す。

【0053】アドレス変換用テーブルは、論理グループテーブル21とアドレステーブル31の2個のテーブルにより構成される。

【0054】図4は、論理グループテーブル21の構造を示す。論理アドレス(論理Addr)22は、CPU1から指定されるアドレスであるところの論理アドレス(データ名またはデータ番号)である。サブデータ名24は、対応する論理アドレス22のデータを分割して得たサブデータの名称である。パリティ名25は、対応するサブデータから作成されたパリティデータの名称である。論理グループ番号(No.)23は、これらのサブデータおよびパリティデータが実際に格納されているまたは格納する論理グループ9の番号である。

【0055】データを全く書き込んでいない初期状態において、論理グループ番号(No.)23の欄にはあらかじめ設定されている論理グループの番号が初期値として設定されている。初期状態では、論理アドレス22、サブデータ名24、およびパリティ名25は、すべて空欄で空きを示すことになる。

【0056】図5は、アドレステーブル31の構造を示す。アドレステーブル31は、論理グループのどのドライブのどの位置にデータが実際に書き込まれているか(または書き込むか)、その詳細なアドレス情報を保持する。

【0057】図5のアドレステーブル31において、論理アドレス(論理Addr)22は、図4の論理グループテーブル21の論理アドレス22と同じであり、CP

U1から指定されるアドレスであるところの論理アドレス（データ名またはデータ番号）である。また、論理グループ番号（No.）23も図4の論理グループテーブル21の論理グループ番号23と同じである。30は、対応する論理アドレスのデータが実際に書き込まれている（または書き込む）SCSIドライブのアドレス情報を示す。

【0058】SCSIドライブアドレス情報30は、4列のアドレス情報からなる。各列は、ボードアドレス（Addr）27、ドライブ番号（Drive No）28、障害フラグ29およびボード抜去フラグ100からなる。この内、ボード抜去フラグ100は、そのボードアドレスのボードがマザーボードから抜き去られているか否かを示す。32は、実際にデータが各ドライブのどこに書き込まれているか（または書き込むか）を示すドライブ内アドレスである。論理アドレス22に対応するボードアドレス27、ドライブ番号28、およびドライブ内アドレス32により、その論理アドレス22のデータ（サブデータとパリティデータ）が書き込まれている（または書き込む）位置が、どのボードのどのドライブのどの位置かを知ることができる。

【0059】各ボードに対するボード抜去フラグ100は、そのボードをマザーボード37から引き抜いたときに1にセットされる。ユーザは何れかのボードを引き抜いたとき、ユーザがCPU1を介して、このボードを引き抜いたことを示すコマンドをアレイコントローラ2に発行するようになっている。MP1.7はこのコマンドに応答して、そのコマンドが示すボードに対応するボード抜去フラグ100を1にセットする。また、後に、この引き抜かれたボードが、マザーボードに再度挿入されたときには、同様にしてユーザがCPU1からのコマンドにより、このボードが搭載されたことをアレイコントローラ2に通知する。アレイコントローラ2は、このときに、後に述べるような回復処理をした後に、このボードに対するボード抜去フラグ100をリセットする。

【0060】障害フラグ29は対応するボード内の対応するドライブに障害が発生したため、読み出しまたは書き込みができない場合オン（1）、正常の場合オフ（0）となり、ドライブの障害の有無を判定することが可能となる。

【0061】SCSIドライブアドレス情報30の4つの列は、図4の論理グループテーブル21のサブデータ名24およびパリティ名25の列に順に対応している。

【0062】例えば、図4の論理グループテーブル21によれば、論理アドレスD#1に対応する論理グループ番号23はLG#1、サブデータ名24とパリティ名25は順にD#1-1、D#1-2、D#1-3、P#1である。したがって、論理アドレスD#1のデータは、3つのサブデータD#1-1、D#1-2、D#1-3に分割されてパリティデータP#1が作成され、それら

が論理グループLG#1の各ドライブに格納されていることが分かる。

【0063】さらに、図5のアドレステーブル31により、各サブデータおよびパリティデータが実際に格納されている位置が分かる。具体的には、サブデータD#1-1はボードアドレス27がBADR1のボードのドライブ番号28がSD#1のドライブ5に格納され、サブデータD#1-2はボードアドレス27がBADR2のボードのドライブ番号28がSD#1のドライブ5に格納され、サブデータD#1-3はボードアドレス27がBADR3のボードのドライブ番号28がSD#1のドライブ5に格納され、これらのサブデータにより作成されたパリティデータP#1はボードアドレス27がBADR4のボードのドライブ番号28がSD#1のドライブ5に格納されていることが分かる。さらに、これらのサブデータおよびパリティは、各々のドライブ5内のドライブ内アドレス32がSADR1の位置に格納されていることが分かる。

【0064】データを全く書き込んでいない初期状態において、図5のアドレステーブル31の論理グループ番号23の欄には、あらかじめ設定されている論理グループの番号が初期値として設定されている。また、初期状態において、ボードアドレス（Addr）27、ドライブ番号（Drive No）28、およびドライブ内アドレス32の欄もあらかじめ設定されている。初期状態では、論理アドレス22がすべて空欄となっており、該当する位置が空き領域であることになる。

【0065】上記のアドレス変換用テーブル（論理グループテーブル21とアドレステーブル31）は、動作時には、アレイコントローラ内のキャッシュメモリ4内の適当な領域に格納されている。これらのアドレス変換用テーブルは、システムの電源をオンしたときに、MP1.7により論理グループ9内のある特定のドライブ5からキャッシュメモリ4に自動的に読み込まれる。一方、電源をオフするときは、MP1.7によりキャッシュメモリ4内のアドレス変換用テーブルが、元のドライブ5内の所定の場所に自動的に格納される。このような電源オン/オフの際のアドレス変換用テーブルのロードとストアはCPU1の関与なく行なわれる。

【0066】なお、本実施例では、同じドライブ内アドレス32を有する複数の領域にデータを並列に格納するため、論理グループ9内の各ドライブ5の回転をすべて同期させる方が望ましい。

【0067】本実施例でのデータの書き込みと読み出しは、このようなアドレステーブル31、および論理グループテーブル21を使用して行なわれる。

【0068】（4）ドライブに障害が発生したときの装置動作

本実施例によれば、図6からも分かるように、複数のボード（4枚のボードBADR1、BADR2、BADR

3, B A D R 4) にまたがる複数のドライブで論理グループ9を構成するようにしており、図4, 5のアドレス変換用テーブルにおいては各ドライブのアドレスにボードアドレスを付加して、どのボードのどのドライブで論理グループを構成しているのが明確になるようにしている。

【0069】本実施例では、何れかのドライブ、例えばB A D R 1上のドライブに読み出しまたは書き込み処理を行なおうとしたが、処理が行なえず、何回かリトライを行なっても処理できない場合、M P 1 7は、ボードB A D R 1に障害が発生したと判断し、アドレステーブル31のボードB A D R 1のドライブS D # 1に対応する障害フラグ29をオン(1)にする(図5)。このように、このドライブに障害が発生したとき、そのドライブを搭載するボードB A D R 1がマザーボード37に搭載した状態でも、このボードをマザーボード37から取り外して修理に供している状態でも、このディスクアレイシステムは引き続き使用できるようになっている。

【0070】すなわち、C P U 1からこの各ボードB A D R 1-4上の各ドライブS D # 1に保持されたデータに対して読み出し要求が発行された場合、M P 1 7は、アドレステーブル31を調べ、アドレス変換用テーブルに関連して先に述べたような手順で、C P U 1から指定された論理アドレスをディスクアレイシステム内のアドレスに変換する。このとき、その変換後のアドレスが割り当てられているドライブ、今の例ではボードB A D R 1上のドライブS D # 1に対応する障害フラグ29と、このドライブが搭載されているボードに対応するボード抜去フラグ100を調べる。このドライブS D # 1に対して障害フラグ29がオン(1)で、このドライブが搭載されているボードB A D R 1のボード抜去フラグ100がオフ(0)の場合は、M P 1 7は、このドライブに障害が発生しており、しかも、このドライブが搭載されているボードB A D R 1はマザーボード37に接続されていることを認識する。

【0071】このようにドライブの障害をM P 1 7が認識すると、M P 1 7は、このドライブが属する論理グループのデータとパリティデータとから、ボードB A D R 1上のドライブS D # 1のデータを回復する。すなわち、他のボードB A D R 2-4上の他の複数のドライブS D # 1から複数のサブデータおよびパリティデータを読み出し、これらのサブデータおよびパリティデータから障害ドライブS D # 1のサブデータを、M P 1 7はパリティ生成回路18を使用して回復し、先に読み出した複数のサブデータと結合して、要求されたデータを生成し、その生成したデータをC P U 1に送る。障害ドライブS D # 1内のデータの回復をパリティ生成回路18により行なわせることは、それ自体は公知である。

【0072】このように、パリティ生成回路(P G) 18は、データをドライブへ書き込む場合は、サブデータか

らパリティを生成する回路であるが、サブデータとパリティからサブデータを回復する場合はデータ回復回路として使用される。

【0073】一方、ドライブS D # 1の搭載されているボードB A D R 1のボード抜去フラグ100がオン

(1)の場合は、M P 1 7はボードB A D R 1がマザーボード37から抜かれていると認識する。このようにM P 1 7が認識すると、M P 1 7は、このボードB A D R 1に搭載されている全ドライブは、障害がなくても読み出し処理ができない。このため、先に述べた障害が発生したドライブS D # 1での処理と同様にして、それらのドライブがアクセスされたときに、そのアクセスされたドライブ内のデータを、そのドライブと同じ論理グループに属する他のドライブのデータから回復する。

【0074】また、C P U 1からドライブS D # 1に保持されたデータを更新する書き込み要求がきた場合、M P 1 7は、読み出し要求での処理と同様に、アドレステーブルを調べ、アドレス変換用テーブルのところで述べたような手順で、C P U 1から指定された論理アドレスをディスクアレイシステム内のアドレスに変換する。このとき、該当するドライブに対応する障害フラグ29と、この当該ドライブが搭載されているボードに対応するボード抜去フラグ100を調べる。

【0075】ドライブS D # 1に対して障害フラグ29がオン(1)で、このドライブS D # 1が搭載されているボードB A D R 1のボード抜去フラグ100がオフ

(0)の場合は、M P 1 7は、ドライブS D # 1に障害が発生しており、しかも、このドライブS D # 1が搭載されているボードB A D R 1はマザーボード37に接続されていると認識する。このように、M P 1 7が認識すると、M P 1 7は、この書き込みデータはキャッシュメモリ4内に保持しておく。

【0076】後に、このボードB A D R 1の障害が発生したドライブS D # 1の修理が完了した時点で、このボードB A D R 1がマザーボード37に再度組み込まれたときに、この障害が発生したドライブに代わる正常なドライブ、今の例ではボードB A D R 1の新たなドライブS D # 1と同じ論理グループに属する他のドライブのすべての誤り訂正データ群に属するサブデータおよびパリティデータを読み出し、それらのサブデータおよびパリティデータからその障害が発生したドライブ内の全データを回復し、この障害が発生したドライブを置換した正常な交替ドライブに格納する。

【0077】その後、キャッシュメモリ4に保持されていた複数の書き込みデータに対して、書き込み動作を実行して、それぞれの書き込むデータを、上記ボードB A D R 1上の、上記交替ドライブあるいは他の障害が発生しなかった正常なドライブに書き込む。この書き込みを行なうときには、R A I D 3による書き込みを行なうために、この保持された各書き込みデータを複数のサブデ



ータに分割し、それらからパリティデータを生成し、分割により得られた複数のサブデータの一つを上記ボード上の一つのドライブに書き込み、他の複数のサブデータと生成されたパリティデータを他のボード上の複数のドライブに書き込む。

【0078】一方、ドライブSD#1が搭載されているボードBADR1のボード抜去フラグ100がオン

(1)の場合は、MP17はボードBADR1がマザーボード37からぬかれていると認識する。このようにMP17が認識すると、MP17は、このボードBADR1に搭載されている全ドライブは、障害の有無にかかわらず書き込み処理ができないため、先に述べた障害が発生したドライブSD#1での処理と同様にドライブSD#1に書き込むデータをキャッシュメモリ4内に保持しておく。

【0079】後に、この障害が発生したドライブを正常な交替ドライブで置換した後のドライブボードBADR1がマザーボード37に再度組み込まれたときに、この障害が発生したドライブを置換した交替ドライブと同じ論理グループに属する他のドライブのすべての誤り訂正データ群に属するサブデータおよびパリティデータを読み出し、それらのサブデータおよびパリティデータからその障害が発生したドライブ内の全データを回復して、その交替ドライブに書き込む。その後、キャッシュメモリ4に保持されていた書き込みデータに対して、書き込み動作を実行する。

【0080】なお、本実施例では、分離されたボードBADR1内の障害が発生していないドライブも、このボードが分離されている間は、アクセスできなくなる。これらのドライブに対するデータの書き込みあるいは読み出しも、その障害ドライブを搭載したボードBADR1がマザーボード37から分離されている場合に、その障害があるドライブに対して行なうのと同じように行なう必要がある。ただし、この分離されたボード上の正常なドライブに対して書き込み要求が発生した場合には、そのボードが後に再度マザーボードに接続された時点では、上記障害が発生したドライブを置換したドライブに対して行なうと説明した、そのドライブ内の全データの回復を、その正常なドライブに関して行なう必要はない。そのような正常なドライブは、そのボードを分離するまでのデータを正常に保持しているからである。

【0081】このように新たに発生した書き込みデータをキャッシュメモリ4に保持するので、マザーボード37から分離されていた間に発生した書き込みデータを、分離されたボード内の正常なドライブに書き込むことが容易に行ない得る。

【0082】また、キャッシュメモリ4に書き込みデータを保持しているので、後にその保持されたデータを読み出す要求がCPU1から供給された場合、キャッシュメモリ4からその保持されたデータをCPU1に読み出

しデータとして供給すればよいので、読み出しが実質的に高速になるという利点も有する。

【0083】以上のごとくにして、本実施例では、障害発生により何れかのドライブボードが取り去られた後でも、そのドライブに関連するデータの読み出しおよび書き込み要求を引き続き処理することができる。

【0084】(実施例1の変形例1) 障害があるドライブを保持するボードを分離した状態で、CPU1から供給された、そのボード上のドライブ障害があるドライブSD#1に対する書き込み要求(あるいはそのボード上の他の正常なドライブに対する書き込み要求)は、以下のように処理することも可能である。

【0085】すなわち、この書き込みデータをサブデータに分割し、それからパリティを生成するところまでは、ボードBADR1が分離されていない場合と同様に行なう。さらに、ボードBADR1が分離されている場合でも、このボード上のドライブSD#1(あるいは他の正常なドライブ)が属する論理グループの内、この障害が発生したドライブ(あるいは正常なドライブ)を有するボードBADR1以外のボードBADR2-4内のドライブSD#1あるいは他のドライブに、それぞれのサブデータあるいは新たに生成されたパリティを書き込む。分離されたボードBADR1内のドライブSD#1あるいは他の正常なドライブに書き込むべきサブデータは捨てる。

【0086】後に、このボードBADR1が、マザーボード37に再度接続された時点で、実施例1で述べたように、この分離されたボードに保持されたすべてのドライブに対して、それぞれのドライブのデータの回復を行なう。

【0087】この回復により、障害が発生したドライブ内の元のデータおよびそのボードを分離中に発生した書き込み要求に付随する書き込みデータを、このボードの、この障害ドライブを交替する新たなドライブに回復するとともに、このボードの他の正常なドライブにも、障害発生時点の前からそれらの正常なドライブが保持していたデータおよびそのボードの分離中に発生した新たな書き込み要求に付随する書き込みデータを回復できる。

【0088】上記方法では分離されていたボードBADR1がマザーボード37に再接続された際、ボードBADR1上の正常なドライブに対し全データの回復を行なう。この方法では、ボードBADR1が分離中に書き込み要求が発生していないデータも回復する。ボードBADR1が分離中に非常に多くの書き込み要求が発生し、正常なドライブ内の大部分のデータが更新されている場合は有効である。しかし、分離中に書き込み要求があまり発生しなかった正常なドライブに対しては、すでに保持されている書き込みデータばかりを回復することになり、回復処理時間が無駄である。

【0089】そこで、分離中に書き込み要求があまり発生しなかった正常なドライブに対しては、すでに保持されていた書き込みデータは回復せず、分離中に上位装置から発行された書き込みデータのみを回復する。

【0090】この変形例では、分離されたボード上のドライブの数が多きときには、それらのドライブのデータの全体を回復するのに時間がかかるが、実施例1のように、新たな書き込みデータをキャッシュメモリ4に保持しなくてもよく、それだけキャッシュメモリ4のサイズが小さくてもよいという利点がある。

【0091】（実施例1の変形例2）変形例1と同様に、障害が発生したドライブを有するボードBADR1をマザーボードから分離した状態において、CPU1から書き込み要求が供給されたときに、そのデータを複数のサブデータに分割し、これらの複数のサブデータからパリティデータを生成する。これらの複数のサブデータの内、上記ボードBADR1内のドライブに書き込むべきサブデータ以外と、上記生成されたパリティデータを、上記ボードBADR1以外の、互いに異なる複数のボード上の複数のドライブに書き込む。実施例1の変形例1と異なり、これらのサブデータのの一つが上記ボードBADR1上の、障害が発生していない正常なドライブに書き込むべきサブデータであるときには、そのサブデータをキャッシュメモリ4に保持する。後にボードBADR1が再度マザーボードに接続された時点で、実施例1の変形例1と同様に、ボードBADR1上の、障害ドライブを置換した交替ドライブが保持すべき全データを他のドライブのデータを利用して回復する。一方、ボードBADR1に保持された正常ドライブには、実施例1と同様に、キャッシュメモリ4に保持されたサブデータを書き込む。

【0092】この方法によれば、キャッシュメモリに保持されるサブデータの量が実施例1より少なくなり、かつ、ボードBADR1がマザーボードに再度接続された時点で、ボードBADR1上の正常なドライブに対して、実施例1の変形例1で行なったデータの回復は行なう必要がない。

【0093】（実施例1の変形例3）本実施例1では、バス38を複数のドライブボードで共通で利用される高速なバスとしたが、各ドライブボードとアレイコントローラ2とを1対1で結ぶ、そのドライブボード専用のバスとしてもよい。このように専用のバスにすると、複数のドライブ5を同時に並列に動作させることが可能となる。

【0094】（実施例2）本実施例では、図1あるいは図3に点線で示した予備ボード50Aがさらに付加されている点で実施例1と異なる。図2(b)では、500Aはこの予備ボード用のコネクタを例示する。予備ボード50Aの構造は他のドライブボードと同じであり、そのボード上に他のボード上にあるドライブと同数のドラ

イブが搭載され、予備ドライブとして使用される。本実施例では、各予備ドライブは、他の複数のボード上に搭載されたドライブにより形成されている複数の論理グループの一つに対応し、その対応する論理グループ内の何れかのドライブに障害が発生したとき、その障害が発生したドライブの代わりに使用される。以下では、この予備ボード以外のドライブボードの番号を、実施例1と同様に、BADR<sub>i</sub> (i=1, 2, 3または4) で表し、この予備ボードのボード番号をBADR5とする。

【0095】この予備ドライブの存在に伴い、本実施例では論理グループテーブル21として図7に示すものが、図3に示すものに代えて使用される。本実施例では、予備ドライブに格納するデータ名の欄26が設けられている。アドレステーブル31には、図5に代えて図8に示すものが使用される。このアドレステーブル31は、各論理グループを構成するドライブの総数は、一つの予備ドライブを含む、総計5つのドライブからなる点で図5と異なる。

【0096】本実施例における、ドライブに障害が発生したときの動作は以下の通りとなる。以下では、ボードBADR1上のドライブSD#1に障害が発生したときの装置動作を説明する。

【0097】(A) 障害ドライブに関連する装置動作  
(A1) 障害ドライブを有するボードがマザーボードに接続されている場合

(A1a) 障害ドライブに保持されたデータの読み出し  
障害ドライブを有するボードBADR1がマザーボード37に接続された状態で、CPU1からこの障害ドライブに保持されたサブデータを一部に含むデータに対して読み出し要求が発生されたとき、実施例1の場合と同様に、この障害ドライブと同じ論理グループ内に属する他のドライブのサブデータおよびパリティから障害ドライブ内のサブデータを回復する。今の例ではボードBADR2内のドライブSD#1、ボードBADR3内のドライブSD#1から二つのサブデータおよびボードBADR4内のドライブSD#1内のパリティデータとから、ボードBADR1の障害ドライブSD#1に保持されたサブデータを回復する。その回復されたサブデータを、他の読み出された二つのサブデータと結合して、CPU1に送る。

【0098】本実施例では、この後の動作が実施例1と異なる。すなわち、この回復されたサブデータを、この論理グループに対して設けられた予備ドライブ、今の例ではボードBADR5のドライブSD#1に書き込む。その後、このサブデータに関しては、この予備のドライブを障害の発生したドライブに代えて使用する。すなわち、この後、上の読み出し要求で指定されたデータに対して、CPU1から再度読み出し要求が発行されたときには、あるいはこのデータに対して書き換え要求が出されたとき、この予備のドライブを含む4つの正常なドラ

イブに対して、読み出し要求あるいは書き込み要求を実行する。

【0099】以上の動作のために、予備ボードBADR5内の予備ドライブSD#1内のアドレスSADR1の位置を、障害ドライブの代わりに使用するように、アドレステーブル31Aを書き換えるのは言うまでもない。

【0100】以上の動作から分かるように、この予備ドライブを使用することにより、一度アクセスされたデータをその後再度アクセスするときには、実施例1で必要であったサブデータの回復が不要になる。さらにそのデータを書き換えるときには、実施例1のごとく、書き込みデータをキャッシュメモリ4に保持しておく必要がなく、その書き込み要求を実行できる。

【0101】(A1b) 障害ドライブに保持されたデータに対する書き込み動作

障害ドライブを有するボードがマザーボード37に搭載された状態で、この障害ドライブが保持するサブデータを含むデータに対して、CPU1からデータ書き込み要求が発行された場合、実施例1と異なり、この書き込み要求を即実行する。

【0102】すなわち、通常の書き込み動作と同様に、この書き込みデータを分割して3つのサブデータを得、それらのサブデータからパリティデータを生成する。

【0103】通常の書き込み動作と異なるのは、これらのサブデータの内、障害が発生した、ボードBADR1内のドライブSD#1に書き込むべきサブデータは、予備ボードBADR5内のドライブSD#1に書き込むことである。

【0104】もちろん、この動作のために、アドレステーブル31Aを予備ドライブを使用するように書き換える。

【0105】(A2) 障害ドライブを有するボードがマザーボードから分離接続されている状態での装置動作

この場合の障害ドライブに保持されたデータの読み出しおよびデータ書き込み動作は、障害ドライブがマザーボードに接続されている場合の、データ読み出し動作(A1a)およびデータ書き込み動作(A1b)と同じである。

【0106】(A3) 障害ドライブを有するボードをマザーボードに復元した後の装置動作

障害ドライブが正常なドライブにより置換された後、障害ドライブを有するボードBADR1がマザーボード37に再度搭載された後では、実施例1と同様に、元の障害ドライブが保持していたデータを回復し、障害ドライブを置換したドライブに書き込む。この回復に当たっては、予備ドライブ内に保持されているデータを使用する必要はない。なお、この方法によれば、障害が発生した後に、新たにCPU1から書き込まれたデータも回復される。

【0107】したがって、この障害ドライブの代わりに

使用した、障害ドライブを有するボードBADR5上の予備ドライブは、この後は、同じ論理グループに属する他のドライブに障害が発生したときにその新たな障害ドライブの代わりに使用されることになる。

05 【0108】(B) 障害ドライブを有するボード上の正常ドライブに関連する装置動作

(B1) 障害ドライブを有するボードがマザーボードに接続されている場合

実施例1と同様に、障害ドライブを有するボード上の正常ドライブはそのまま引き続き使用される。

10 【0109】(B2) 障害ドライブを有するボードがマザーボードから分離されている場合

この場合には、障害ドライブを有するボード上の正常なドライブはアクセスできない。したがって、前述の障害ドライブの動作(A1)で説明された動作を行なう。

15 【0110】(B3) 障害ドライブを有するボードをマザーボードに復元した後の装置動作

この場合、分離されていた障害ドライブを有するボード上の正常なドライブが保持していた全データの内、障害ドライブを有するボードがマザーボードから分離した後に、この正常ドライブに書き込み要求が発生している場合、その正常なドライブの元のデータは使用できない。この書き込み要求に伴って生成された、この正常なドライブに書き込むべきデータは、前述した通り、この正常なドライブの予備のドライブに書き込まれている。したがって、障害ドライブを有するボードがマザーボード37に再度接続された後は、この正常なドライブに対する予備ドライブ内に保持されている有効なデータを、障害ドライブを有する、再度接続されたボード内の正常なドライブに移動すればよい。こうして、分離された障害ドライブを有するボード上の正常なドライブに対するデータの回復は完了する。

25 【0111】なお、障害ドライブを有するボードが分離されている間に、この正常なドライブに保持されているデータに対して読み出し要求がCPU1より発行された場合、すでに述べたごとく、この正常なドライブに保持されているデータが回復され、このドライブに対する予備ドライブに格納されている。したがって、上記データの移動により、このような読み出し要求に伴い回復されたデータも、予備ドライブから正常なドライブに移動されることになる。しかし、このように読み出し要求によって回復されたデータは、もともと正常なドライブに保持されているデータと同じである。したがって、上記移動のときに、このようなデータを移動しても、移動完了後の、正常なドライブのデータは正常な値を有する。

40 【0112】このように、本実施例では、実施例1と異なり、予備のドライブが設けられ、何れかのドライブボードがマザーボード37から分離されている間、CPU1から供給された読み出し要求に応答して回復されたデータが、この予備のドライブに保持される。したがって

て、そのようなデータがその後CPU1から要求されたときに、すぐに利用可能になるという利点を有する。もちろん、キャッシュメモリ4に書き込みデータを保持する実施例1と異なり、本実施例は、書き込みデータを予備のドライブに保持するので、より多くの書き込みデータを保持できるという利点がある。

【0113】(実施例2の変形例) 実施例2に対して、障害ドライブを有するボードをマザーボードに復元した後のデータの回復に関連して、次のようないろいろの変形が可能である。

【0114】(1) 障害ドライブのデータの回復に関連する変形

(1a) 上述の実施例2では、障害ドライブを置換したドライブが保持すべき全データを、それが属する論理グループ内の他のドライブ内の複数のサブデータとパリティデータとを使用して回復し、その置換したドライブに格納した。

【0115】障害ドライブの予備のドライブには、障害ドライブが保持すべきデータの内の一部のデータがすでに保持されている。したがって、実施例2における、障害ドライブを置換したドライブへのデータの回復に当たり、この予備のドライブに保持されているデータ以外のデータは、それらを回復して、その置換したドライブに格納するが、予備のドライブにすでに保持されているデータは、その予備のドライブからその置換したドライブに移動すれば、実施例2より高速にデータの回復を完了できる。この方法は、予備のドライブに格納されているデータ量が多いときに有効である。

【0116】(1b) 実施例2では、障害ドライブを有するボードが復帰したとき、もとの障害ドライブを置換したドライブに、もとの障害ドライブが保持すべきであったデータを回復し、その障害ドライブに対して使用した予備ドライブは、新たに予備ドライブとして解放した。しかし、この方法に代えて、障害ドライブが保持すべきデータを回復して、予備のドライブに書き込み、この予備のドライブを引き続き正常なドライブとして使用することもできる。この際、もとの障害ドライブが保持すべき全データの内、この予備ドライブにまだ保持されていないデータを選択的に回復して、この予備のドライブに書き込むことが回復時間を短縮する上で望ましい。

【0117】この方法は、もとの障害ドライブが保持すべきデータの内、かなり大きな割合のデータがすでに予備ドライブに保持されている場合に有効である。

【0118】(1c) 実施例2では、上位装置から障害ドライブ内のデータに対する読み出し要求が供給されたときに、その供給を契機として、そのデータを回復し、さらに、この回復後のデータを予備ドライブに格納した。しかし、この予備ドライブへの格納を省略できる。すなわち、予備ドライブには、障害があるドライブを搭載したボードがマザーボードから分離されている間に上

位装置から供給された書き込みデータだけを保持する。もちろん、このボードがマザーボードに再度接続された後は、障害ドライブのデータを回復して予備ドライブに格納し、その予備ドライブをその障害ドライブの代わりに使用してもよい。

【0119】(1d) 何れかのドライブに障害が発生したときに、直ちに上位装置からディスクアレイへのアクセスを禁止し、その障害ドライブの全データを回復し、予備ドライブに格納し、その後は、この予備ドライブを障害が発生したドライブの代わりに使用してもよい。

【0120】(2) 障害ドライブを有するボード上の正常ドライブに関連する変形例

(2a) 実施例2では、障害ドライブを有するボードをマザーボードに復元した後、障害ドライブを有するボード上の正常なドライブに対する予備のドライブにあるデータをすべて、その正常なドライブに移動した。しかし、この予備のドライブにあるデータには、障害ドライブを有するボードが分離されている間にCPU1から発生された、その正常なドライブのデータに対する読み出し要求を契機として回復されたデータも含まれている。このようなデータはすでにその正常なドライブに保持されている、その障害ドライブを有するボードを分離される前のデータと同じ値を有する。したがって、実施例2における、予備ドライブから、正常ドライブへのデータの移動のときに、このような読み出しを契機として回復されたデータは移動しないようにすることが、移動を速く完了する上で望ましい。

【0121】このためには、予備のドライブにデータが書き込まれたとき、そのデータが書き込み要求を契機として書き込まれたデータか読み出し要求を契機として書き込まれたデータか否かを示す情報をアドレステーブルに格納し、その情報にしたがって、そのデータを予備ドライブから復帰されたボード上の正常ドライブに移動するか否かを制御すればよい。この方法は、予備ドライブに保持されたデータの内、CPU1からの書き込み要求を契機として書き込まれたデータが、CPUからの読み出し要求を契機として書き込まれたデータより多いときに有効である。

【0122】(2b) 実施例2では、障害ドライブを有するボードが復帰したとき、このボード上の正常なドライブに、そのドライブの予備のドライブに書き込まれたデータを移動した。しかし、この方法に代えて、この正常な復帰したドライブが保持すべきデータを回復して、この予備のドライブに書き込み、この予備のドライブを引き続き正常なドライブとして使用することもできる。この際、この正常なドライブが保持すべき全データの内、この予備ドライブにまだ保持されていないデータを選択的に回復して、この予備のドライブに書き込むことが有効である。

【0123】この方法は、正常なドライブに保持すべき

データの内のかなりの割合のデータがすでに予備ドライブに保持されている場合に有効である。

【0124】（実施例3）実施例1、2はRAID3の制御を適用したディスクアレイシステムであるが、本実施例3は本願発明を実施例1の装置の制御にRAID5を採用した実施例である。したがって、本実施例は、図1に示すものと同じ装置構成により実現されるが、MP17による制御が実施例1と異なる。したがって、本実施例でも図1などの回路を使用して本実施例を説明する。説明は実施例1と異なる点を主に示す。図9は、この実施例で使用するアドレステーブルの一例を示す。

【0125】（ドライブ障害がないときの装置動作）このときの装置の動作は、基本的には、RAID5の動作として周知の動作と同じであるが、RAID5の動作の理解のための簡単にこのときの動作を説明する。

【0126】RAID3では、一つの書き込み要求に付随した書き込みデータが複数のサブデータに分割され、それらのサブデータからパリティデータを生成する。生成された複数のサブデータとパリティデータを互いに異なるドライブに書き込む。RAID5ではRAID3とは異なり、一つの書き込み要求に付随するデータは分割されない。複数の書き込み要求に付随する複数のデータに対してパリティデータが定義される。これらの複数のデータとそのパリティデータが一つのパリティグループを形成する。これらのデータとパリティデータとは互いに異なるドライブに書き込まれる。しかし、同一のパリティグループに属する複数のデータが揃ってからそれらのデータとパリティデータが複数のドライブに書き込まれるのではない。各書き込み要求がCPUから供給されるごとに、そのデータが属するパリティグループに対してすでに生成されたパリティデータ（旧パリティデータ）とその書き込みデータとから新パリティデータを生成する。その書き込みデータと生成されたパリティデータが二つのドライブに書き込まれる。読み出し要求がCPUから供給されたときには、その読み出し要求で指定されるデータを何れか一つのドライブから読み出し、CPU1に供給する。

【0127】すなわち、互いに異なる複数のドライブに属する複数の記憶位置を何れかのパリティグループ用にあらかじめ決めておき、何れかの書き込み要求がCPU1からアレイディスクコントローラ2に供給された時点で、その書き込み要求が指定する論理アドレスに対して、一つのパリティグループと、そのパリティグループに対して割り当てられた何れかのドライブの何れかの記憶領域を割り当てる。そのパリティグループに割り当てられた他の記憶領域に未だ何れのデータも保持されていないときには、そのパリティグループの他のデータは所定の値、例えば0とみなし、その書き込みデータからパリティデータを生成する。上記書き込みデータとこの生成されたパリティデータを互いに異なるドライブに書き

込む。

【0128】次に来た第2の書き込み要求に対して同じパリティグループが割り当てられたとき、そのパリティグループに対して生成されたパリティデータ（旧パリティデータ）を読み出し、先にその第2の書き込み要求が指定する第2のデータとこの旧パリティから新たにパリティを生成し、この第2の書き込みデータをそのデータに割り当てられた記憶位置に書き込むとともに、生成されたパリティでもって、旧パリティを書き直す。その後も同様にそのパリティグループの全データに対する書き込みを実行する。このようにRAID5では、RAID3と異なり、データの書き込みごとに、そのデータが属するパリティグループの旧パリティデータとその書き込みデータから新パリティデータを生成するが、このパリティの生成にそのグループの他のデータは使用しない。書き込み済みのデータを読み出すときには、そのデータを何れか一つのドライブから読み出すが、そのデータが属するパリティグループの他のデータは読み出さない。

【0129】（ドライブ障害があるときの装置動作）このときの動作はRAID3とRAID5との差に由来する違いを除いて実施例1と基本的には同じである。例えば、障害ドライブを有するボードがマザーボードに接続されている状態では、CPU1からの、正常なドライブへのデータの書き込み要求を実行する場合には、上記ドライブ障害がない場合に説明したのと同様に、この書き込みデータを分割する必要はない。また、この障害ドライブを有するボードがマザーボードから分離されている間に、この分離されているボード上の障害ドライブに対し上位装置から書き込み要求が発行された場合、この書き込み要求に付随する書き込みデータは一時的にキャッシュメモリ4に保持しておく。

【0130】その障害ドライブが正常な交替ドライブと置換された後、そのボードがマザーボードと再度接続されたときに、交替ドライブが保持すべきデータを回復し、回復されたデータをその交替ドライブに書き込み、その後、キャッシュメモリに保持された書き込みデータを、その書き込みデータを分割することなくその交替ドライブに書き込む。障害ドライブを有するボード内の正常なドライブに関しても、このデータ回復を行なう必要がないことを除いて、同様である。

【0131】（実施例3の変形例）実施例1の各変形例がこの実施例3にも適用できる。ただし、RAID3とRAID5との違いに由来する変更が必要である。

【0132】例えば、実施例1の変形例1では、障害ドライブを有するボードがマザーボードから分離されている状態で障害ドライブに対する書き込み要求がCPU1から供給された場合、書き込みデータを複数のサブデータに分割し、それからパリティデータを生成し、これらのサブデータの内、障害ドライブに書き込むべきデータ

以外のサブデータおよび生成された新パリティデータをそれぞれ異なるドライブに書き込んだ。実施例 1 のこの変形例 1 に対応する本実施例 3 の変形例では、書き込みデータを分割する必要はない。すなわち、障害ドライブを有するボードがマザーボードから分離されている間に上位装置から書き込みデータが供給されたとき、その書き込みデータが属するパリティグループの複数のデータの内、その障害があるドライブ以外のドライブに保持された複数のデータおよびそのパリティグループに属するパリティデータを読み出し、これらの読み出されたデータと新たに CPU 1 から供給された書き込みデータとから、新パリティを生成する。生成された新パリティデータでもって旧パリティを書き換えるだけでよい。

【0133】また、RAID 4 のアレイコントローラにも本実施例は適用できる。

【0134】（実施例 4）本実施例は本願発明を実施例 2 の装置の制御に RAID 5 を採用した実施例である。したがって、本実施例は、実施例 2 と同様に図 1 に示すものに予備ドライブボード 50A を付加したものと同一装置構成により実現されるが、MP 17 による制御が実施例 2 と異なる。したがって、本実施例の説明でも図 1 などの回路を使用して説明する。なお、説明は実施例 2、3 と異なる点を主にする。

【0135】（ドライブ障害がないときの装置動作）このときの動作は RAID 3 と RAID 5 との差に由来する違いを除いて実施例 2 と基本的には同じである。その違いは、実施例 3 と実施例 1 との違いに関して説明したとおりである。

【0136】（ドライブ障害があるときの装置動作）このときの動作は RAID 3 と RAID 5 との差に由来する違いを除いて実施例 2 と基本的には同じである。その違いは、実施例 3 と実施例 1 との違いに関して説明したとおりである。

【0137】（実施例 4 の変形例）実施例 2 の各変形例がこの実施例 4 の変形例としても使用できる。ただし、RAID 3 と RAID 5 との違いに由来する変更が必要である。

【0138】例えば、障害ドライブへのデータの書き込みに関する CPU 1 からの書き込み要求を予備ドライブを使用して実行する場合には、上記ドライブ障害がない場合と同様に、この書き込みデータを分割する必要がない。すなわち、その書き込みデータが属するパリティグループの複数のデータの内、その障害があるドライブ以外のドライブに保持された複数のデータおよびそのパリティグループに属するパリティデータを読み出し、これらの読み出されたデータと新たに CPU 1 から供給された書き込みデータとから、その書き込みデータに対応する新パリティデータを生成し、この書き込みデータを予備ドライブに書き込み、この新パリティデータでもって旧パリティデータを書き換える。

【0139】また、RAID 4 のアレイコントローラにも本実施例は適用できる。

【0140】（実施例 5）以上の実施例と異なり、本実施例はドライブに障害がないがボードのコネクタに障害が発生した場合の実施例である。

【0141】何れかのドライブボードのコネクタがマザーボードとの接触不良などの不良が発生した場合にも、そのボード上の全ドライブはアクセスできなくなる。その場合、このボードを正常なものと交換するまで、そのボードをマザーボードから抜き取ることが必要になる。本実施例では、このような場合に、以上の実施例のいずれかを適用したものである。

【0142】すなわち、ユーザは何れかのボードのコネクタの接触異常を発見したとき、CPU 1 を介してアレイコントローラ 2 に、そのボードを抜き取った通知するコマンドを入力する。その後、そのボードを抜き取る。この後の装置の動作は、上記実施例の何れかに記載のとおりである。

【0143】なお、何れかのドライブボードのコネクタの接触不良は次のような場合に検出可能である。そのドライブボード内のドライブをアクセスした結果、そのアクセスが不成功に終わると、アレイコントローラ 2 は、そのボードがマザーボードから分離されていない状態では、そのドライブに障害が発生したことを示す障害フラグをアドレステーブルにセットする。このようにして同じボードのドライブが次々と障害が発生したと判断された場合、一般には、それらのドライブの障害よりも別のところの障害の可能性が高い。その障害の原因の一つはそのボードのコネクタの接触不良である。ユーザによる調査の結果、そのコネクタの接触不良が判明したときには、ユーザは上記のコマンドを入力する。それに先立ち、すでに障害が発生したとされたドライブがあれば、それらのドライブに対する障害ビットをオフにするコマンドをユーザが入力する。アレイコントローラ 2 は、このコマンドの入力に回答して、そのコマンドで指定されるボード上の障害とされたドライブ上のドライブに対する障害ビットをリセットする。

【0144】

【発明の効果】以上説明したように、本発明によるディスクアレイ装置は、障害が発生したドライブを正常なドライブに交換する際、障害が発生したドライブが搭載されている基板をマザーボードから外して交換している間も使用し続けることができる。また、何れかのドライブボードのコネクタに障害が発生しても、そのボード上のドライブ内のデータを指定するデータ読み出し要求およびデータ書き込み要求を実行できる。

【図面の簡単な説明】

【図 1】本発明に係るディスクアレイシステムの第 1 の実施例の概略的な装置構成を示す図

【図 2】図 1 の装置のボードの配置およびマザーボード



のコネクタを示す図

【図 3】 図 1 の装置のより詳細な回路図

【図 4】 図 3 の回路に使用する論理グループテーブル (21) の一例を示す図

【図 5】 図 3 の回路に使用するアドレステーブル (31) の一例を示す図

【図 6】 第 1 の実施例によるデータの書き込みを説明する図

【図 7】 本発明に係るディスクアレイシステムの第 2 の実施例で使用する論理グループテーブルの一例を示す図

【図 8】 第 2 の実施例で使用するアドレステーブルの一例を示す図

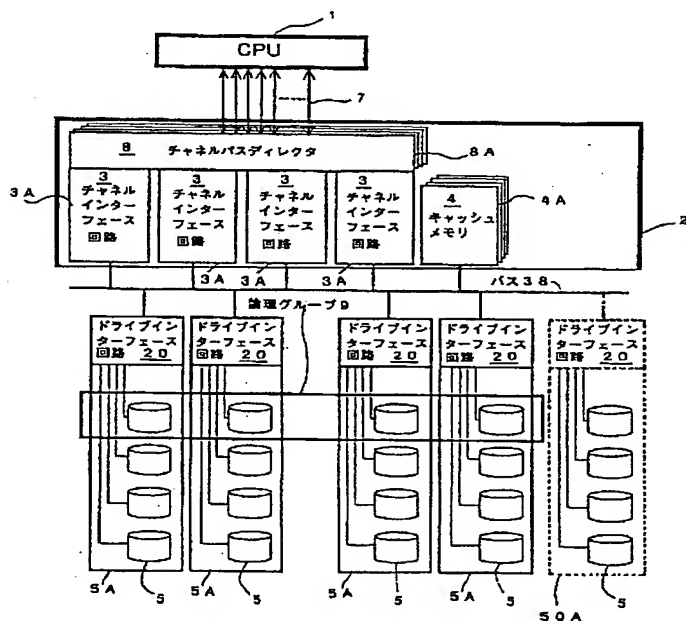
【図 9】 本発明に係るディスクアレイシステムの第 3 の実施例で使用するアドレステーブルの一例を示す図

【符号の説明】

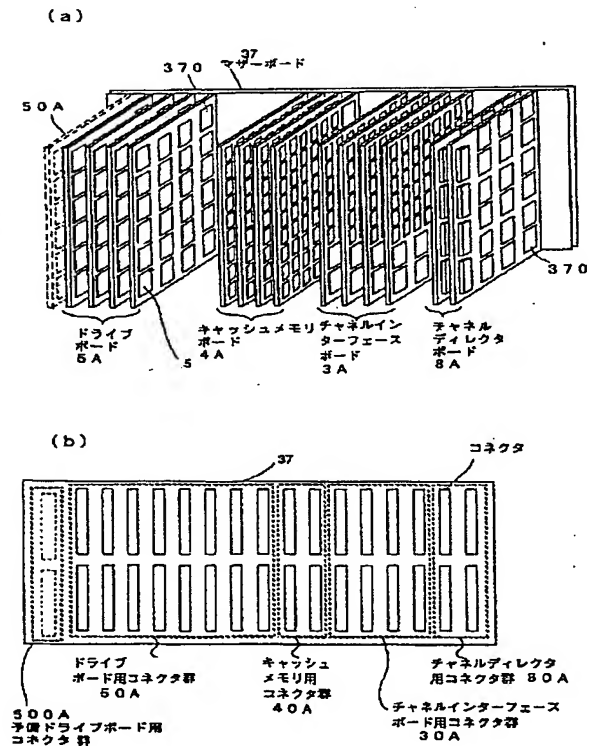
2…アレイコントローラ、3A…チャンネルインターフェースボード (CIB)、4A…キャッシュメモリボー

ド、5A…ドライブボード、7…チャンネルバス、8…チャンネルバスディレクタ、9…論理グループ、10…インターフェースアダプタ、11…チャンネルバススイッチ、12…データ線、13…チャンネルインターフェース (C H I F) 回路、14…データ制御回路 (DCC)、15…チャンネル側キャッシュアダプタ (C Adp)、16…制御信号線、17…マイクロプロセッサ (MP)、18…パリティ生成回路 (PG)、19…ドライブ側キャッシュアダプタ (C Adp)、20…ドライブインターフェース回路 (Drive IF)、21…論理グループテーブル、22…論理アドレス (論理Addr)、23…論理グループNo、24…サブデータ名、25…パリティ名、27…ボードアドレス (ボードAddr)、28…Drive No、29…障害フラグ、30…SCSIドライブアドレス、31…アドレステーブル、32…Drive内アドレス、100…ボード抜きフラグ。

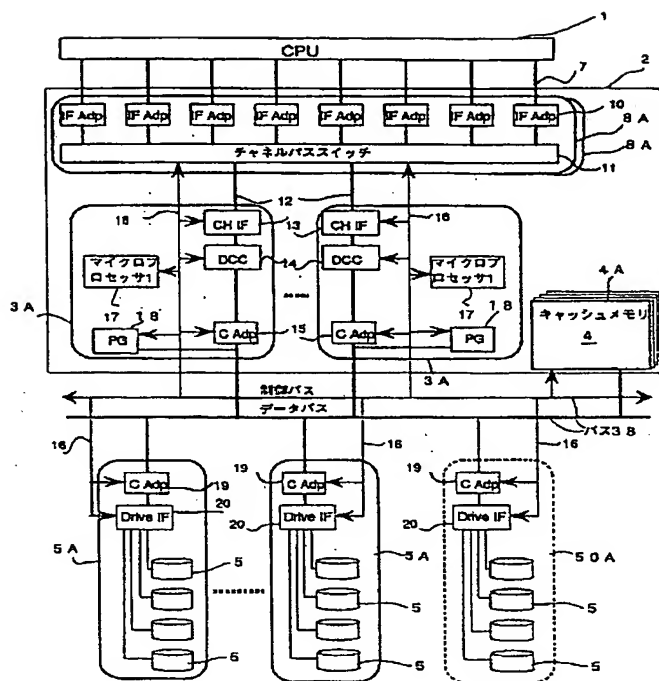
【図 1】



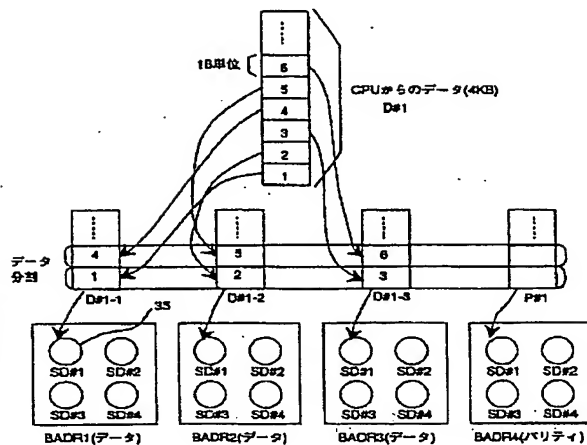
【図 2】



【図3】



【図6】



【図4】

論理グループテーブル

22 論理 Addr	23 論理グル ープNo.	21 24 サブデータ名			25 パリティ 名
D#1	LG#1	D#1-1	D#1-2	D#1-3	P#1
D#2		D#2-1	D#2-2	D#2-3	P#2
D#3		D#3-1	D#3-2	D#3-3	P#3
D#4		D#4-1	D#4-2	D#4-3	P#4
D#5		D#5-1	D#5-2	D#5-3	P#5
⋮		⋮	⋮	⋮	⋮
D#11	LG#2	D#11-1	D#11-2	D#11-3	P#11
D#12		D#12-1	D#12-2	D#12-3	P#12
D#13		D#13-1	D#13-2	D#13-3	P#13
D#14		D#14-1	D#14-2	D#14-3	P#14
D#15		D#15-1	D#15-2	D#15-3	P#15
⋮		⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮

【図 8】

100

31 アドレステーブル

31 アドレステーブル

[illegible]

【図7】

論理グループテーブル

22 論理 Addr	23 論理グル ープNo.	21 サブデータ名			24 パリティ 名	25 Spare
D#1	LG#1	D#1-1	D#1-2	D#1-3	P#1	Spare1
D#2		D#2-1	D#2-2	D#2-3	P#2	Spare2
D#3		D#3-1	D#3-2	D#3-3	P#3	Spare3
D#4		D#4-1	D#4-2	D#4-3	P#4	Spare4
D#5		D#5-1	D#5-2	D#5-3	P#5	Spare5
⋮		⋮	⋮	⋮	⋮	⋮
D#11	LG#2	D#11-1	D#11-2	D#11-3	P#11	Spare11
D#12		D#12-1	D#12-2	D#12-3	P#12	Spare12
D#13		D#13-1	D#13-2	D#13-3	P#13	Spare13
D#14		D#14-1	D#14-2	D#14-3	P#14	Spare14
D#15		D#15-1	D#15-2	D#15-3	P#15	Spare15
⋮		⋮	⋮	⋮	⋮	⋮
⋮		⋮	⋮	⋮	⋮	⋮

【図9】

3.1 アドレステーブル

論理グループ No.	論理 Addr	SCSI ドライブアドレス							
		データ				パリティ			
		ボード Addr	ボード No.	Drive No.	ボード No.	ボード Addr	ボード No.	Drive No.	Drive Addr
LG#1	D#1	BADR1	0	SD#1	0	BADR4	0	SD#1	0
	D#2	BADR2	0		0	BADR4	0	SD#1	0
	D#3	BADR3	0		0	BADR4	0	SD#1	0
	D#4	BADR2	0		0	BADR4	0	SD#1	0
	D#5	BADR3	0	SD#1	0	BADR1	0	SD#1	0
	D#6	BADR4	0		0	BADR1	0	SD#1	0
	D#7	BADR3	0		0	BADR1	0	SD#1	0
	D#8	BADR4	0		0	BADR2	0	SD#1	0
	D#9	BADR1	0	SD#1	0	BADR2	0	SD#1	0
	D#10	BADR4	0		0	BADR2	0	SD#1	0
	D#11	BADR1	0		0	BADR3	0	SD#1	0
	D#12	BADR2	0		0	BADR3	0	SD#1	0
	D#13	BADR1	0	SD#2	0	BADR4	0	SD#2	0
	D#14	BADR2	0		0	BADR4	0	SD#2	0
	D#15	BADR3	0		0	BADR4	0	SD#2	0
	D#16	BADR2	0		0	BADR1	0	SD#2	0
	D#17	BADR3	0	SD#2	0	BADR1	0	SD#2	0
	D#18	BADR4	0		0	BADR1	0	SD#2	0
	D#19	BADR3	0		0	BADR2	0	SD#2	0
	D#20	BADR4	0		0	BADR2	0	SD#2	0
	D#21	BADR1	0	SD#2	0	BADR3	0	SD#2	0
	D#22	BADR4	0		0	BADR3	0	SD#2	0
	D#23	BADR1	0		0	BADR3	0	SD#2	0
	D#24	BADR2	0		0	BADR3	0	SD#2	0
	:	:	:	:	:	:	:	:	:
	:	:	:	:	:	:	:	:	: